

AD-A112 043

CALIFORNIA UNIV BERKELEY DEPT OF CIVIL ENGINEERING F/6 12/1  
A NEWTON-LANCZOS METHOD FOR SOLUTION OF NONLINEAR FINITE ELEMEN--ETC(U)  
FEB 82 B NOUR-OMID N62583-82-MR-419

UNCLASSIFIED

CEL-CR-82.012

NL

100  
100

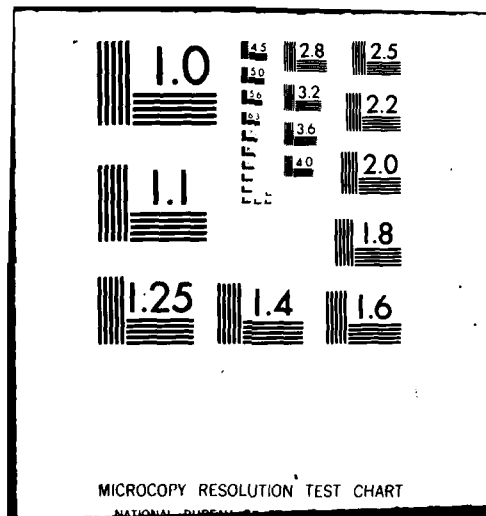
END

DATE

FILED

4-82

DTIC



ADA112043

Unclassified

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER CR 82.012	2. GOVT ACCESSION NO. AD-A112 043	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) A Newton-Lanczos Method for Solution of Nonlinear Finite Element Equations		5. TYPE OF REPORT & PERIOD COVERED Not Final; Mar 1981 - Oct 1981
7. AUTHOR(s) Bahram Nour-Omid		8. PERFORMING ORG. REPORT NUMBER
9. PERFORMING ORGANIZATION NAME AND ADDRESS Department of Civil Engineering University of California, Berkeley, CA 94720		10. CONTRACT OR GRANT NUMBER(s) N62583-82-MR-419
11. CONTROLLING OFFICE NAME AND ADDRESS Naval Civil Engineering Laboratory Port Hueneme, CA 93043		12. PROGRAM ELEMENT PROJECT, TASK AREA & WORK UNIT NUMBERS YR023.03.01.003
13. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		14. REPORT DATE February 1982
		15. NUMBER OF PAGES 63
		16. SECURITY CLASS (of this report) Unclassified
		17. DECLASSIFICATION DOWNGRADING SCHEDULE
18. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlimited		
19. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)		
20. SUPPLEMENTARY NOTES		
21. KEY WORDS (Continue on reverse side if necessary and identify by block number) Structural mechanics, Finite elements, Nonlinear algebraic equations, Numerical solution methods		
22. ABSTRACT (Continue on reverse side if necessary and identify by block number) The finite element method reduces nonlinear continuum problems to nonlinear discrete problems which take the form of systems of nonlinear algebraic equations. Attention is de- voted to procedures which may be employed to solve the result- ing nonlinear algebraic systems. The general class of continu- um problems of interest include both material and geometric nonlinearities.		

DTIC  
ELECTE  
MAR 16 1982  
H

Robert A

Unclassified

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

→ Newton's method, modified Newton methods, and quasi-Newton methods are reviewed. However, the technique which has been given focus is the Newton-Lanczos method which is a member of a class of solution methods that employ an iterative, linear equation solver in an inner loop within Newton's method.

The Newton-Lanczos algorithm is shown to not only require fewer factorization steps than either the quasi-Newton or modified Newton methods but also possesses more robust convergence characteristics when dealing with nearly singular Jacobian matrices and indefinite systems.

Unclassified

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

### ACKNOWLEDGEMENTS

This research work was sponsored by Naval Civil Engineering Laboratory under Contract No. N62583/81 M R419.

I would like to express my gratitude to Professor R. L. Taylor who supervised this work and gave me encouragement and freedom to pursue my own related interests. He provided support and much valuable advice on this project. I am deeply indebted to Professor B. N. Parlett for giving me the privilege of working under his supervision. I appreciate his patience, understanding and generosity with his valuable time during the course of this research. I am also grateful to Professors J. L. Sackman and A. Der Kiureghian for their encouragement and many helpful and enjoyable discussions.



Accession For	
NTIS GRA&I	<input checked="" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	
By	
Distribution/	
Availability Codes	
Avail and/or	
Dist	Special
A	

## TABLE OF CONTENTS

<b>ACKNOWLEDGEMENTS .....</b>	<b>i</b>
<b>1. INTRODUCTION .....</b>	<b>1</b>
<b>2. THE FINITE ELEMENT METHOD FOR NONLINEAR PROBLEMS .....</b>	<b>3</b>
2.1 Equations of Kinematics .....	3
2.2 Momentum Balance .....	4
2.3 Constitutive Relation .....	5
2.4 Finite Element Discretization .....	6
<b>3. SOLUTION OF THE DISCRETE NONLINEAR EQUATIONS .....</b>	<b>10</b>
3.1 Newton's Method .....	10
3.1.1 Line Search for Newton's Method .....	12
3.1.2 Operation Counts for Newton's Method .....	12
3.1.3 Convergence of Newton's Method .....	13
3.1.4 Advantages and Disadvantages of Newton's Method .....	13
3.2 Modified Newton's Method .....	14
3.3 Quasi-Newton Methods .....	15
3.3.1 Broyden's Method .....	16
3.3.2 Computation of the Inverse Matrix .....	17
3.3.3 Convergence of Quasi-Newton Method .....	17
3.3.4 Methods for Symmetric Positive Definite Jacobians .....	18
<b>4. LANCZOS ALGORITHM FOR SOLUTION OF THE     LINEARIZED PROBLEM .....</b>	<b>20</b>
4.1 The Lanczos Algorithm .....	21

4.2	The Weak Form .....	23
4.3	Termination Criteria .....	24
4.4	Derivation of Conjugate Gradients Method from Lanczos Algorithm .....	28
4.5	Selective Orthogonalization .....	29
4.6	Lanczos with Selective Orthogonalization .....	32
4.7	Selective Reorthogonalization .....	33
4.8	Starting Vector .....	34
5.	NEWTON-LANCZOS METHOD .....	39
5.1	Description of the Algorithm .....	39
5.2	Convergence of Newton-Lanczos .....	41
5.3	Preconditioning of The Linear Equations .....	43
5.4	Numerical Results .....	44
5.5	Closure .....	46
APPENDIX A .....		55
APPENDIX B .....		57
REFERENCES .....		58



## 1. INTRODUCTION

During the past few years, considerable effort has been directed towards the solution of nonlinear initial-boundary value problems in structural mechanics. The finite element method is one of the more popular approaches employed to reduce continuum problems to nonlinear algebraic, discrete problems. Finite element methodology is well documented in the literature (e.g., see[23]) and attention here is devoted to procedures which may be employed to solve the resulting nonlinear algebraic problem. The general class of problems of interest include both material and geometric nonlinearities.

A brief review of the application of the finite element to problems in nonlinear continuum mechanics is presented in Section 2. The result is a large set of nonlinear algebraic equations which must be solved for the state variables (e.g., displacements and stresses). Newton's method is commonly employed to solve the nonlinear equations. In Section 3 discussion of Newton's method, including operation count estimates, and advantages and disadvantages, is presented. In addition modified Newton's method is also discussed in this section. In general, Newton's method possesses ideal characteristics of local convergence and stability but is normally too expensive to employ for solving large finite element problems. On the other hand, modified Newton's method has desirable operation count estimates for each iteration but becomes expensive to employ because of its low convergence rate characteristics. Recently, Matthies and Strang [10] have suggested that quasi-Newton methods be used to solve nonlinear finite element problems. In Section 3.3 some of the quasi-Newton's methods which have been used in optimization methods, such as Broyden's method and the Broyden-Fletcher-Goldfarb-Shanno (BFGS) method, are discussed and a discussion of the advantages and disadvantages for these methods is presented.

The use of iterative equation solvers as an inner loop for the Newton's method has been suggested before and recently Eisenstat completed a study of the convergence properties of these methods. In section 4 the use of the Lanczos algorithm for iterative solution of

linear system of equations is described. The relation of the Lanczos algorithm to the conjugate gradient method and the effect of preconditioning on the algorithm is demonstrated. The preconditioned Lanczos algorithm is used to solve the linear system of equations arising at each iterate of Newton's method. This results in a technique where the rate of convergence of the method can be varied from a linear to a quadratic convergence rate by means of a specified tolerance. In section 5 a detailed description of the Newton-Lanczos method is presented. The above methods are tested on various nonlinear problems in structural mechanics and the results are tabulated.

## 2. THE FINITE ELEMENT METHOD FOR NONLINEAR PROBLEMS

Since its introduction, the Finite Element Method, FEM, has become one of the most widely used techniques for solving various problems in continuum mechanics. The FEM has established a systematic procedure for discretization of complex structures which results in a large set of nonlinear equations.

The source of nonlinearities are of two types:

- (a) material nonlinearity
- (b) geometrical nonlinearity

Addressed here are problems associated with both nonlinear constitutive relations and nonlinearities due to large displacements.

In the following section a summary of the formulation of problems in continuum mechanics is presented. The resulting equations are then discretized using the Galerkin method. The usual convention, that implies summation with respect to a repeated index in a term, is adopted.

### 2.1. Equations of Kinematics

Consider a body with reference configuration at time  $t = 0$ , denoted by  $B_0$ , that is deformed by a motion to a current configuration  $B$ . Figure 2.1 illustrates this motion. Using a common rectangular cartesian coordinate system, the Lagrangian description of the motion can be expressed as

$$\mathbf{x} = \hat{\mathbf{x}}(\mathbf{X}, t) = \mathbf{X} + \mathbf{u}(\mathbf{X}, t) \quad (2.1.1)$$

where  $\mathbf{x}$  is the current position at time  $t$  of a particle with position vector  $\mathbf{X}$  at time  $t = 0$  and  $\mathbf{u}$  is the displacement the particle undergoes in this motion.

The deformation gradient of the motion is defined by

$$\mathbf{F} = \frac{\partial \mathbf{x}}{\partial \mathbf{X}} = \mathbf{I} + \frac{\partial \mathbf{u}}{\partial \mathbf{X}} \quad (2.1.2)$$

The matrix  $\frac{\partial \mathbf{u}}{\partial \mathbf{X}}$  is known as the displacement gradient. Thus, the strain tensor, at a point in the current configuration, referred back to the initial configuration can be written as

$$\mathbf{E} = \frac{1}{2}(\mathbf{F}^T \mathbf{F} - \mathbf{I}) \quad (2.1.3)$$

Using equations 2.1.2 and 2.1.3 the strain tensor can be written directly in terms of the displacement and the initial position vector

$$\mathbf{E} = \frac{1}{2} \left[ \left( \frac{\partial \mathbf{u}}{\partial \mathbf{X}} \right)^T + \frac{\partial \mathbf{u}}{\partial \mathbf{X}} + \left( \frac{\partial \mathbf{u}}{\partial \mathbf{X}} \right)^T \frac{\partial \mathbf{u}}{\partial \mathbf{X}} \right] \quad (2.1.4.a)$$

or in indicial notation

$$E_{IJ} = \frac{1}{2}(u_{I,J} + u_{J,I} + u_{K,I}u_{K,J}) \quad (2.1.4.b)$$

where  $E_{IJ}$  are components of the strain tensor  $\mathbf{E}$ , and  $u_I$  are components of displacement vector  $\mathbf{u}$ ;  $(\ )_{,J}$  denotes partial differentiation with respect to the  $J$ -th component of the position vector  $\mathbf{X}$ .

For transient problems a statement of initial conditions is required. That is

$$\mathbf{u}(\mathbf{X}, 0) = \mathbf{d}_0(\mathbf{X}) \quad (2.1.5.a)$$

$$\dot{\mathbf{u}}(\mathbf{X}, 0) = \mathbf{v}_0(\mathbf{X}) \quad (2.1.5.b)$$

## 2.2. Momentum Balance

Equilibrium conditions between internal resisting forces and externally applied forces must be satisfied, at least in the weak form, whether the displacements or strains are large or small. The momentum balance equation is given by

$$\text{Div} \mathbf{P} + \mathbf{b} = \rho_0 \ddot{\mathbf{u}} \quad (2.2.1.a)$$

or

$$P_{Ii,I} + b_i = \rho_0 \ddot{u}_i \quad (2.2.1.b)$$

where  $P_{ji}$  are the components of the first Piola-Kirchhoff (unsymmetric) stress tensor  $\mathbf{P}$ ,  $b_i$  are components of the body force  $\mathbf{b}$ ,  $\rho_0$  is the density in the initial configuration, and a superposed dot,  $(\dot{\phantom{x}})$ , denotes differentiation with respect to time.

A weak form of the momentum balance equations, equivalent to virtual work, may be constructed by multiplying 2.2.1.b by an arbitrary function  $W_i$ , integrating over the domain of interest  $\Omega$ , using integration by parts on the stress term and setting the result equal to zero.

The result is

$$\int_{\Omega} (W_{i,j} P_{ji} - W_i b_i + W_i \rho_0 \ddot{u}_i) d\Omega - \int_{\Gamma_2} W_i \bar{\tau}_i d\Gamma = 0 \quad (2.2.2)$$

where in addition to previously defined quantities,  $\bar{\tau}_i$  is a specified traction and  $\Gamma_2$  is the part of the boundary where traction is specified. For equation 2.1.2 to be valid,  $W_i$  must vanish on  $\Gamma_1$ , that part of the boundary where displacements are specified (i.e.,  $u_i = \bar{u}_i$  on  $\Gamma_1$ ). The traction is related to stress through

$$\bar{\tau}_i = N_j P_{ji} \quad (2.2.3)$$

where  $N_j$  are direction cosines of the outward normal to  $\Gamma$  in the initial configuration.

The components of the first Piola-Kirchhoff stress tensor,  $P_{ji}$ , is replaced by  $S_{ji}$ , the components of the symmetric second Piola-Kirchhoff stress tensor,  $\mathbf{S}$ , using the relation

$$P_{ji} = S_{ji} F_{ji} \quad (2.2.4)$$

where  $F_{ji}$  are the components of the deformation gradient matrix. This results in

$$\int_{\Omega} (W_{i,j} S_{ji} F_{ji} - W_i b_i + W_i \rho_0 \ddot{u}_i) d\Omega - \int_{\Gamma_2} W_i \bar{\tau}_i d\Gamma = 0 \quad (2.2.5)$$

### 2.3. Constitutive Relation

In general the deformation gradient of the motion, the temperature and the temperature gradient determines the stress state. The effects of inelastic material behavior can be modeled through internal variables.

For isothermal deformation the constitutive equation can be written as

$$S_{IJ} = \rho_0 \frac{\partial \hat{\Psi}}{\partial E_{IJ}} \quad (2.3.1)$$

where  $\Psi$  is the Helmholtz free energy defined at  $(X, t)$  through the strain tensor and the internal variables

$$\Psi(X, t) = \hat{\Psi}(E, \alpha) \quad (2.3.2)$$

The form of the constitutive relation stated above is currently under debate. A more comprehensive discussion can be found in reference [14]. In addition to the previously defined quantities  $\alpha$  is a vector of internal variables. Further, a constitutive relation of the form

$$\dot{\alpha} = \gamma(S, \alpha) \quad (2.3.3)$$

is assumed for the internal variables. It is sometimes more convenient to express these relations with strain rate as the dependent variable

$$\dot{E} = C\dot{S} + \Lambda\dot{\alpha} \quad (2.3.4)$$

where  $C$  is the instantaneous elastic compliance tensor, and  $\Lambda$  is the inelastic compliance tensor which depends on the stress state, the internal variables and the velocity gradient,  $L$ , in such a way as to satisfy objectivity

$$\Lambda = \hat{\Lambda}(S, \alpha, L) \quad (2.3.5)$$

This form of the model is most suitable for nonlinear viscoelastic or elastic/viscoplastic materials. The constitutive equations are both nonlinear and rate dependent. Writing 2.3.4 in indicial notation

$$\dot{E}_{IJ} = C_{IJKL} \dot{S}_{KL} + \Lambda_{IJ(R)} \dot{\alpha}_{(R)} \quad (2.3.6)$$

where  $R=1, 2, \dots, N$  and  $N$  is the number of internal variables.

The weak form of this equation can similarly be obtained by multiplying 2.3.6 by an arbitrary function  $V_{IJ}$ , and integrating over the domain  $\Omega$ . This leads to

$$\int_{\Omega} V_{IJ}(\dot{E}_{IJ} - C_{IJKL}\dot{S}_{KL} - \Lambda_{IJ(R)}\dot{\alpha}_{(R)}) d\Omega = 0 \quad (2.3.7)$$

where  $\dot{\alpha}_{(R)}$  is determined from equation 2.3.3 .

#### 2.4. Finite Element Discretization

In a finite element model, the domain  $\Omega$  is divided into elements  $\Omega_e$ . By sampling the solution for the field variables, such as displacements or stresses, at a number of points in the elements, known as nodes, an approximate solution can be constructed. When the stresses are known explicitly in terms of the strains it is possible to choose the displacements as the only independent variable and the solution to the nodal parameter can be obtained using the weak form of the equilibrium equation, 2.2.5 . However, in cases when the constitutive relation is more complex it is simpler to use a weak form of this equation, 2.3.7 , and use both stresses and displacements as primary variables.

In each element the approximations for  $u_i$  and  $S_{IJ}$  which are  $C^0$ -continuous and piecewise constant (  $C^{-1}$ -continuous ), respectively, over  $\Omega_e$ , are in the form

$$u = \sum_A N_A(X) u_A(t) \quad (2.4.1.a)$$

and

$$S = \sum_B M_B(X) S_B(t) \quad (2.4.1.b)$$

where  $N_A(X)$  and  $M_B(X)$  are shape functions over the domain  $\Omega_e$  satisfying the  $C^0$  and  $C^{-1}$  continuity requirements cited above. The arbitrary weighting functions are also approximated using these shape functions, as

$$W = \sum_A N_A(X) W_A \quad (2.4.2.a)$$

and

$$V = \sum_B M_B(X) V_B \quad (2.4.2.b)$$

where  $W_A$  and  $V_B$  are arbitrary nodal parameter. Using a temporal finite difference scheme, such as Newmark's method [13], the time derivatives in equations 2.2.5 and 2.3.7 can be removed.

The result of the application of the finite element method and the time-stepping procedure is a large set of nonlinear algebraic equations at each time step,  $t_m$ , with nodal displacements and stress quantities as the unknowns. The dependence of these equations on the stress quantities can be removed through static condensation at the element level resulting in nodal displacements as the unknown parameters.

These equations can be written as

$$f[x(t_m)] = 0 \quad (2.4.3)$$

where  $x$  is a vector of the state variables, which consists of all the nodal displacements,  $u_j(t_m)$ , and  $f$  is a vector function of  $x$  obtained by application of the Galerkin method to equations 2.2.5 and 2.3.7.



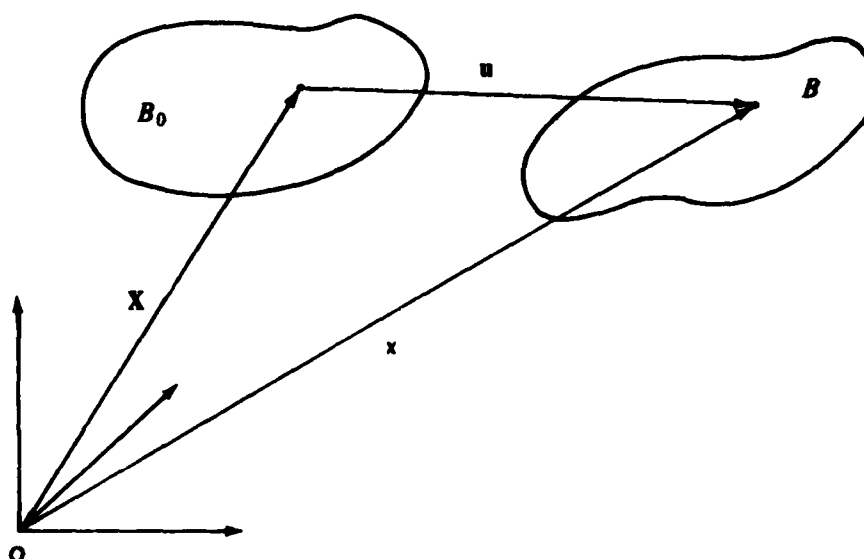


Figure 2.1 Composition of Motion

### 3. SOLUTION OF THE DISCRETE NONLINEAR EQUATIONS

Employing the finite element method and a time stepping procedure to discretize the nonlinear structural mechanics problem leads to the large set of nonlinear algebraic equation given by 2.4.3 . Presented here, are a number of methods normally employed for the solution of this equation.

The Taylor expansion for this , truncated to include only the linear terms can be written as

$$f(x + \Delta x) \approx f(x) + f'(x)\Delta x \quad (3.0.1)$$

where  $f' = \frac{\partial f}{\partial x}$  is a matrix with coefficients  $f'_{ij}$  given by

$$f'_{ij} = \frac{\partial f_i}{\partial x_j}$$

and  $\Delta x$  is a small change in  $x$ .

Almost all nonlinear procedures are based on this approximation.

#### 3.1. Newton's Method

Using equation 3.0.1 Newton's method can be written as

$$f(x_{k+1}) \approx f(x_k) + A_k d_k \quad (3.1.1)$$

where  $k$  is an iteration number,  $d_k$  is a change in  $x_k$  called the step direction, and  $A_k$  is the Jacobian or tangent matrix of  $f$  defined by

$$A_k = \left[ \frac{\partial f(x)}{\partial x} \right]_{x_k} \quad (3.1.2)$$

Newton's method consists of setting 3.1.1 equal to zero, solving for  $d_k$  from

$$A_k d_k = -f(x_k) \quad (3.1.2)$$

and setting

$$x_{k+1} = x_k + d_k \quad (3.1.3)$$

Newton's method requires the initial guess  $x_0$  to be in a domain  $D$ , called the domain of attraction, such that convergence to the exact solution,  $x^*$  in  $D$ , of equation 2.4.3 does occur. Furthermore,  $f$  must be differentiable in  $D$  and the tangent matrix must be non-singular at the solution. In practice it is often desirable to modify 3.1.3 to

$$x_{k+1} = x_k + s_k d_k \quad (3.1.4)$$

where  $s_k$  is a scalar step size which is used to enhance stability of the algorithm. The value of  $s_k$  is determined from a line search as described in the following section.

The algorithm for implementing Newton's method consists of choosing a good initial guess  $x_0$  and repeat following steps for  $k = 0, 1, 2, \dots$  until convergence is achieved:

- (1) given  $x_k$  compute  $f(x_k)$
- (2) compute the tangent matrix  $A_k$
- (3) solve  $A_k d_k = -f(x_k)$  for  $d_k$
- (4) compute  $s_k$  from a line search
- (5) update  $x_{k+1} = x_k + s_k d_k$
- (6) test for convergence and terminate if converged

There are several procedures which may be used to terminate the iteration.

These include:

- (a)  $\|f(x_k)\| < \epsilon_1 \max_i \|f(x_i)\|$
- (b)  $\|d_k\| < \epsilon_2 \max_i \|d_i\|$
- (c)  $|d_k \cdot f(x_k)| < \epsilon_3 \max_i |d_i \cdot f(x_i)|$

where  $\epsilon_i$  are small positive constants. In this work method (a) was used. Method (c) is only applicable for symmetric positive definite Jacobians.

### 3.1.1. Line Search for Newton's Method

In the above presentation a step size,  $s_k$ , has been included to improve the stability of Newton's method. The magnitude of  $s_k$  is obtained in such a way as to minimize a norm of the residual,  $f(x_k)$ . A common procedure used, is to solve the following nonlinear scalar equation for  $s_k$

$$G(s_k) = d_k^T f(x_k + s_k d_k) = 0 \quad (3.1.5)$$

It is important to note that equation 3.1.5 need not be solved accurately since  $x_{k+1}$  is itself only an approximation and also, in general, the additional function evaluations are costly in finite element analyses.

In fact, when possible, the line search should be omitted when  $x_k$  is close to a solution.

### 3.1.2. Operation Counts for Newton's Method

For purposes of subsequent cost comparisons, operation counts for Newton's method are estimated to indicate the relative efforts needed in each step. For an  $n$ -dimensional  $x$ , the operation count estimates are:

(a) computation of $f(x_k)$	$O(n)$
(b) computation of $A_k$	$O(n^2)$
(c) direct solution of equations	
triangular decomposition of $A_k$	$O(n^3)$
forward reduction/backsubstitution	$O(n^2)$
(d) line search	$O(n)$
(e) update of solution	$O(n)$
(f) convergence test	$O(n)$

While order of magnitude estimates are given, substantial differences in real effort exist between, for example (a) and (d). The majority of effort is associated with steps (a) to (d), and only comparisons between these steps are meaningful. There is an order of magnitude

increase in step (b) over step (a) or in step (c) over step (b). It is important to note that these estimates are somewhat problem dependent, for example step (a) can cost as much as step (b). Consequently, considerable efficiency can be achieved by eliminating or reducing the number of expensive steps required.

### 3.1.3. Convergence of Newton's Method

The rate of convergence of an algorithm will determine how fast the iteration  $x_k$  approaches a solution  $x^*$ . An acceptable algorithm must be at least linearly convergent [3]; that is, given a solution  $x^*$ , then

$$\|x_{k+1} - x^*\| \leq \alpha \|x_k - x^*\| \quad (3.1.6)$$

where  $\alpha$  is a positive constant less than unity. Although 3.1.6 ensures that the error norm is reduced by the factor  $\alpha$  in each iteration, to be competitive it is generally acknowledged that an algorithm must have better than linear convergence. When Newton's method has continuously differentiable  $f$  and a solution  $x^*$  in  $D$ , then the error norm satisfies the stronger condition

$$\|x_{k+1} - x^*\| \leq \alpha \|x_k - x^*\|^q \quad (3.1.7)$$

where  $q$  is a positive scalar such that  $1 < q < 2$ . This is called super-linear convergence. In addition, for cases where  $f$  is twice differentiable in the neighborhood of  $x^*$ , with non-singular Jacobian, Newton's method is quadratically convergent with  $q = 2$ .

For finite element applications, Newton's method will almost always have at least super-linear convergence, and most problems are such that quadratic convergence will be achieved.

### 3.1.4. Advantages and Disadvantages of Newton's Method

Newton's method has at least two very desirable properties:

- (i) Any  $x_k$  in  $D$  results in an  $x_{k+1}$  in  $D$ ; consequently, the method is stable and convergent once an iterate is in  $D$ .
- (ii) The method possesses at least super-linear convergence and often quadratic convergence [3].

On the negative side, we note that:

- (i) If  $D$  is small, then the initial approximation may have to be very close to  $x^*$  to get inside  $D$ . We do not know when  $x_0$  is in  $D$ .
- (ii) The triangular decomposition of the Jacobian matrix is very costly in large finite element problems.

The requirement of a good initial guess may be avoided in part by using line searches and, for quasi-static problems, an evolution of the load application [23]. In the sequel we address the possibility of reducing computational factorizations of the tangent matrix.

### 3.2. Modified Newton's Method

For large systems of equations, the main cost in Newton's method is the triangular decomposition of the Jacobian matrix. The use of a previously factored Jacobian is often advocated in place of the current tangent matrix. The algorithm is given as, for  $k=0,1,2, \dots$  repeat

- (1) given  $x_k$  compute  $f(x_k)$
- (2) solve  $A_k d_k = -f(x_k)$ , where  $A_k$  is the tangent matrix at step  $i < k$
- (3) compute  $s_k$  from a line search
- (4) update  $x_{k+1} = x_k + s_k d_k$
- (5) test for convergence and terminate if converged

Comparing this algorithm with that of Newton's method, we observe that the  $O(n^2)$  operations to compute the Jacobian matrix and the  $O(n^3)$  operations to compute the triangular

decomposition are avoided. The algorithm still requires  $O(n^2)$  operations to perform a resolution using the triangular factors of  $A_j$ . However, these savings are achieved at the expense of the convergence rate since modified Newton's method only converges linearly, as given by equation 3.1.6. There exists some number of iterations  $p$  such that computation of a new Jacobian matrix will make the modified Newton's method more efficient. Unfortunately, the value of  $p$  depends on the degree of nonlinearity of the problem and cannot be estimated easily.

### 3.3. Quasi-Newton Methods

Quasi-Newton methods are a generalization of the one-dimensional secant method to the  $n$ -dimensional problem (e.g., see [11] and [12]). In the secant method, an approximation to the tangent matrix,  $A_k$ , is used at each iteration. This concept is applied in multi-dimensions and a simple updating is deduced to compute the new  $A_k$  from the previous value of  $A_{k-1}$ . The starting matrix  $A_0$  is normally taken as  $f'(x_0)$ , however, other choices are possible such as the finite difference approximations [3]. The convergence rate for all practical quasi-Newton methods is super-linear, and the number of numerical operations for each iteration is  $O(n^2)$ .

To deduce the equation for the updating of the Jacobian, known as the secant equation, a linear backward Taylor formula at step  $k$  is used.

$$f(x_k) \approx f(x_{k-1}) - f'(x_k)(x_k - x_{k-1}) \quad (3.3.1)$$

This equation can be rearranged to obtain the approximate Jacobian  $A_k$

$$A_k d_{k-1} = b_{k-1} \quad (3.3.2)$$

where

$$d_{k-1} = x_k - x_{k-1} \quad (3.3.3)$$

and

$$b_{k-1} = f(x_k) - f(x_{k-1}) \quad (3.3.4)$$

The values in 3.3.3 and 3.3.4 must be computed to perform the step

$$A_k d_k = -f(x_k) \quad (3.3.5)$$

Equation 3.3.2 is called the quasi-Newton equation, which must be satisfied by  $A_k$ . In a one dimensional problem, this equation is sufficient to determine  $A_k$  completely, however, for multi-dimensional problems additional constraints are required to obtain  $A_k$ . The construction of the approximate Jacobian depends, to a large extent, on the choices for these constraints. The work of Broyden and Broyden-Fletcher-Goldfarb-Shanno, BFGS, are but two of the many possibilities for constructing  $A_k$ .

### 3.3.1. Broyden's Method

In 1965 Broyden proposed a method for the approximate specification of the tangent matrix by a simple update of the previous value. Broyden assumed that  $A_k$  does not differ from  $A_{k-1}$  when acting on any vector orthogonal to  $d_{k-1}$ , accordingly

$$A_k z = A_{k-1} z \quad (3.3.6)$$

for

$$z^T d_{k-1} = 0$$

These equations provide an update relation for  $A_k$  and, together with equation 3.3.2, allows  $A_k$  to be determined uniquely. Thus, Broyden's method is given by

$$A_k = A_{k-1} + \frac{(b_{k-1} - A_{k-1} d_{k-1}) d_{k-1}^T}{\|d_{k-1}\|^2} \quad (3.3.7)$$

Post-multiplying 3.3.7 by  $d_{k-1}$  and  $z$  results in equations 3.3.6 and 3.3.2, respectively, thus demonstrating the applicability of 3.3.7. This equation can be written in the computationally more attractive form, using equation 3.3.4



$$A_k = A_{k-1} + \frac{f(x_k)d_{k-1}^T}{\|d_{k-1}\|} \quad (3.3.8)$$

Broyden's method is super-linearly convergent [3]. Given an initial  $x_0$  and  $A_0$  and using the above equations the quasi-Newton step can be performed. The algorithm is identical to that for Newton's method, except the computation of the Jacobian is replaced by equation 3.3.8. The factorization of  $A_k$  at each step is still required; consequently, it is best to directly update the inverse of  $A_{k-1}$  to obtain the inverse of  $A_k$ . This will eliminate the need to compute the triangular decomposition of  $A_k$  and therefore will result in an algorithm with  $O(n^2)$  operations in each step.

### 3.3.2. Computation of the Inverse Matrix

The inverse of matrices defined similar to 3.3.8 may be written as

$$(A + vw^T)^{-1} = A^{-1} - \frac{1}{\sigma} A^{-1}vw^T A^{-1} \quad (3.3.9)$$

where

$$\sigma = 1 + w^T A^{-1}v$$

If we let  $H_k$  be the inverse of  $A_k$ , Broyden's method for the update of the inverse becomes

$$H_k = H_{k-1} + \frac{(d_{k-1} - H_{k-1}b_{k-1})d_{k-1}^T H_{k-1}}{d_{k-1}^T H_{k-1} b_{k-1}} \quad (3.3.10)$$

provided  $d_{k-1}^T H_{k-1} b_{k-1}$  is nonzero. Broyden's method may be implemented as

$$d_k = -H_k f(x_k) \quad (3.3.11)$$

### 3.3.3. Convergence of Quasi-Newton Method

Convergence properties of quasi-Newton method are discussed by Dennis and More in [3]. In this study they restate an earlier result that an iterative method is super-linearly convergent if

$$\lim_{k \rightarrow \infty} \frac{\| [A_k - f'(x^*)](x_{k+1} - x_k) \|}{\| x_{k+1} - x_k \|} = 0 \quad (3.3.12)$$

If  $A_k$  converges to  $f'(x^*)$ , as for Newton's method where  $A_k$  equals  $f'(x_k)$ , then convergence is always super-linear. However, the important result of 3.3.12 is that when  $A_k$  only converges to  $f'(x^*)$  in the direction  $d_k$ , convergence is also super-linear. Both the Broyden and the BFGS quasi-Newton method satisfy 3.3.12 for continuously differentiable  $f$  and thus are super-linearly convergent.

### 3.3.4. Methods for Symmetric Positive Definite Jacobians

Brodie, Gourlay and Greenstadt have shown that certain rank-one and rank-two corrections to symmetric positive definite matrices may be conveniently expressed in the product form [1]

$$H_k = (I + w_k v_k^T) H_{k-1} (I + v_k w_k^T) \quad (3.3.13)$$

This form has the important property that successive inverse tangent matrices remain symmetric. In [10] the algorithm is related to the BFGS algorithm, which gives

$$w_k = \frac{1}{d_{k-1}^T b_k} d_{k-1} \quad (3.3.14)$$

and

$$v_k = f(x_{k-1}) \left[ 1 - \left( \frac{d_{k-1}^T b_{k-1}}{d_{k-1}^T f(y_k)} \right)^{1/2} \right] - f(x_k) \quad (3.3.15)$$

The computational steps for implementing the BFGS method are very straight-forward and consist of solving 3.3.11 with following steps:

- (1) for each  $k$  compute  $f(x_k)$  and  $z_k = - (I + v_k w_k^T) f(x_k)$
- (2) solve  $H_{k-1}^{-1} u_k = z_k$
- (3) compute  $d_k = (I + w_k v_k^T) u_k$

- (4) if required, compute a line search for  $s_k$
- (5) update  $x_{k+1} = x_k + s_k d_k$
- (6) check convergence

It should be noted that  $H_{k-1}$  may result from a previous BFGS update. Thus, step (1) and (3) may require several updates before the resolve step is performed. Furthermore, this form of the algorithm is strictly limited to positive definite tangent matrices. Indefinite matrices resulting from Lagrange multiplier constraint, such as contact problems, can not be considered by present implementation of BFGS and therefore an alternative implementation is required.

#### 4. LANCZOS ALGORITHM FOR SOLUTION OF THE LINEARIZED PROBLEM

The discretization of nonlinear structural mechanics problems and linearization of the resulting nonlinear algebraic equations for application of a Newton type methods normally will lead to a symmetric system of linear equations. These equations may be written as

$$Ax = b \quad (4.0.1)$$

Presently, most nonlinear procedures employ a direct method such as Gaussian elimination for the solution of this system of equations. In this section an alternative solution technique is presented.

In most applications  $A$  will be sparse, and an elegant way to make use of sparsity is to employ  $A$  solely as a linear operator which computes the product  $Av$  for a given vector  $v$ . This has the added advantage that  $A$  need not be known explicitly but only a means of computing the matrix vector product is required.

For years, iterative methods have been used for the solution of large systems of equations. The Conjugate Gradient method is one such technique introduced by Hestenes and Stiefel [5]. This method, in theory, requires at most  $n$  steps to obtain the exact solution of 4.0.1 and was accepted because of this termination property. In the same year Lanczos published his method of Minimized Iteration which was initially introduced for computing the eigen pairs of a large symmetric matrix. Lanczos and Householder [6] pointed out the intimate connection between the two approaches.

These methods have several attractive features in common. There is no need for  $A$  to have further special properties, such as banded form, no acceleration parameters have to be estimated, and the storage requirements are only a few  $n$ -vectors in addition to the storage needs of  $A$ .

In finite element applications the matrix vector multiplication can be performed at the element level. This will result in considerable saving in storage at the expense of some

additional multiplications.

#### 4.1. The Lanczos Algorithm

For certain applications of the finite element method, especially in nonlinear problems, it is usual to have on hand an initial approximation  $x^a$  to the true solution of 4.0.1. The problem is now to find a correction  $x^c$  to be added to  $x^a$ . Then

$$Ax^c = r_0 \quad (4.1.1)$$

where

$$r_0 = b - Ax^a \quad (4.1.2)$$

Lanczos algorithm may be described very simply as a process of constructing the weak form of equation 4.0.1 from a very special subspace. The subspace under consideration is generated from the set of  $j$  vectors  $(r_0, Ar_0, \dots, A^{j-1}r_0)$ , known to mathematicians as the Krylov subspace [18]. To construct the weak form it would be simpler if an orthonormal set of vectors, say  $(q_1, q_2, \dots, q_j)$ , were available. This can be achieved by applying Gram-Schmidt orthogonalization to the Krylov vectors. Initially, this appears to be an expensive way of obtaining an orthonormal base vectors, however this process can be simplified when the following two facts are used:

- (i) The use of  $Aq_j$  and  $A^j r_0$  for orthogonalization against the previous  $q$  vectors and normalization of the resulting vector, leads to the same vector  $q_{j+1}$ .
- (ii) The vector  $Aq_j$  is orthogonal to  $q_1, q_2, \dots, q_{j-2}$ .

These results are demonstrated in Appendix A. A more complete explanation is presented in [18]. Therefor it is sufficient to orthogonalize  $Aq_j$  against  $q_{j-1}$  and  $q_j$  to obtain the next orthogonal vector.

$$r_j \equiv \beta_{j+1}q_{j+1} = Aq_j - \alpha_j q_j - \beta_j q_{j-1} \quad (4.1.3)$$

where  $\alpha_j = \mathbf{q}_j^T \mathbf{A} \mathbf{q}_j$  and  $\beta_j = \mathbf{q}_{j-1}^T \mathbf{A} \mathbf{q}_j$ . It is important to note that the vectors  $(\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_{j-2})$  are not needed in equation 4.1.3. This defines one step of the Lanczos algorithm. The normalization of  $\mathbf{r}_j$  results in  $\mathbf{q}_{j+1}$ . It is easy to show, by looking at  $\mathbf{q}_{j+1}^T \mathbf{r}_j$ , that  $\beta_{j+1} = \|\mathbf{r}_j\|$ .

The special choice of the base vectors for the subspace has an additional advantage. The projection of  $\mathbf{A}$  onto this subspace is a tridiagonal matrix,  $\mathbf{T}_j$ .

$$\mathbf{T}_j = \mathbf{Q}_j^T \mathbf{A} \mathbf{Q}_j = \begin{bmatrix} \alpha_1 & \beta_2 & & & \\ \beta_2 & \alpha_2 & \beta_3 & & \\ & \beta_3 & & & \\ & & & \ddots & \\ & & & & \alpha_{j-1} & \beta_j \\ & & & & \beta_j & \alpha_j \end{bmatrix} \quad (4.1.4)$$

where the  $\mathbf{q}$  vectors form the columns of the matrix  $\mathbf{Q}_j$ ,  $\mathbf{Q}_j \equiv (\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_j)$ .

This fact was realized soon after Lanczos introduced his method and the algorithm was put to use as a process for the orthogonal transformation of a matrix to tridiagonal form. Despite its additional attractions, the Lanczos process gave way to Givens' method in 1954 and later to Householder's method in 1958.

The relationships that define the Lanczos algorithm can now be summarized in the following three equations.

$$\mathbf{Q}_j^T \mathbf{Q}_j = \mathbf{I}_j \quad (4.1.5.a)$$

$$\mathbf{A} \mathbf{Q}_j - \mathbf{Q}_j \mathbf{T}_j = \mathbf{r}_j \mathbf{e}_j^T \quad (4.1.5.b)$$

$$\mathbf{Q}_j^T \mathbf{r}_j = 0 \quad (4.1.5.c)$$

where  $\mathbf{e}_j$  is the  $j$ -th column of the  $j \times j$  identity matrix  $\mathbf{I}_j$ .

Setting  $\mathbf{q}_0 = 0$  and using  $\mathbf{r}_0$  as the starting vector, the Lanczos algorithm can then be described as

$$\text{Given } \mathbf{r}_0, \text{ set } \beta_1 = \|\mathbf{r}_0\|$$

For  $j = 1, 2, \dots$  repeat

$$(1) \mathbf{q}_j \leftarrow \frac{\mathbf{r}_{j-1}}{\beta_j}$$

$$(2) \mathbf{u}_j \leftarrow \mathbf{A}\mathbf{q}_j$$

$$(3) \mathbf{r}_j \leftarrow \mathbf{u}_j - \beta_j \mathbf{q}_{j-1}$$

$$(4) \alpha_j \leftarrow \mathbf{q}_j^T \mathbf{r}_j$$

$$(5) \mathbf{r}_j \leftarrow \mathbf{r}_j - \alpha_j \mathbf{q}_j$$

$$(6) \beta_{j+1} \leftarrow \|\mathbf{r}_j\|$$

At the end of a simple Lanczos step the newly computed Lanczos vector is written into secondary storage.

#### 4.2. The Weak Form

A weak form of equation 4.0.1 can now be constructed using equation 4.1.5. The approximation to the solution from the subspace can be written as

$$\mathbf{x}_j = \mathbf{Q}_j \mathbf{f}_j \quad (4.2.1)$$

where  $\mathbf{f}_j$  is the vector in the subspace defining the approximate solution. The  $i$ -th coefficient  $\phi_i$  of  $\mathbf{f}_j$  and its corresponding  $\mathbf{q}_i$  are much like the nodal parameter and its associated shape function in the finite element discretization.

Equations 4.1.1 together with equation 4.2.1 results in

$$\mathbf{A}\mathbf{Q}_j \mathbf{f}_j = \mathbf{r}_0 \quad (4.2.2)$$

This is an overdetermined system of equation. However a solution can be obtained using a procedure similar to the method of weighted residuals, and an approximation can be obtained by premultiplying this with  $\mathbf{W}_j^T$ . The weighting vectors  $\mathbf{W}_j$  can be chosen in such manners as to reduce the residual of the quantities of interest.

$$W_j^T A Q_j f_j = W_j^T r_0$$

A Galarkin type method results when the Lanczos vectors are chosen as the weighting vectors.

$$\begin{aligned} Q_j^T A Q_j f_j &= Q_j^T r_0 \\ T_j f_j &= \beta_1 e_1^{(j)} \end{aligned} \quad (4.2.3)$$

Here  $e_1^{(j)}$  is the first column of the identity matrix and the above is derived using the fact that  $r_0 = q_1 \beta_1$ . The above equations possess the orthogonal characteristics of the residual vector peculiar to Galarkin methods, that is the residual vector is kept orthogonal to the trial vectors and therefore is "minimized". In appendix B equation 4.2.3 is also derived from an energy consideration. This orthogonality property can be seen in the form of a monotonic reduction in the potential energy of the system which is a characteristic of both the Lanczos method and conjugate gradient algorithm. In Figure 4.1 this is shown in terms of a simple beam problem.

#### 4.3. Termination Criteria

In the previous section the Lanczos algorithm was viewed as a method for constructing the weak form of the linear set of equations. The form of the algorithm suggests that after every Lanczos step the number of  $q$  vectors increase by one. In exact arithmetic the algorithm can be terminated after  $n$  steps and the tridiagonal matrix contains all the information needed to solve the problem exactly. However, it is often sufficient to stop the iteration at some step  $m < n$ , when the solution is close enough to the exact solution. Using the weak form, equation 4.2.3, the solution vector

$$f_m = [\phi_1, \phi_2, \dots, \phi_m]^T \quad (4.3.1)$$

can be obtained by solving the equations. The approximate solution  $x_m$  can then be constructed directly from equation 4.2.1. The residual vector is



$$\begin{aligned}
 \hat{r}_m &\equiv Ax_m - r_0 \\
 &= AQ_m f_m - r_0
 \end{aligned} \tag{4.3.2}$$

Setting  $j = 1$  in equation 4.1.3 the following can be derived

$$\begin{aligned}
 r_0 &= \beta_1 q_1 \\
 &= Q_m (e_1^{(m)} \beta^1)
 \end{aligned} \tag{4.3.3}$$

This together with 4.3.2 and 4.2.3 results in

$$\begin{aligned}
 \hat{r}_m &= AQ_m f_m - Q_m T_m f_m \\
 &= (AQ_m - Q_m T_m) f_m \\
 &= (r_m e_m^T) f_m \\
 &= r_m \phi_m
 \end{aligned} \tag{4.3.4}$$

Therefore the residual vector is a scalar multiple of the  $r_m$  vector computed at the  $m$ -th Lanczos step. Moreover the residual norm  $\rho_m$  is given by

$$\begin{aligned}
 \rho_m &\equiv \|\hat{r}_m\| \\
 &= |\phi_m| \|r_m\| \\
 &= \beta_{m+1} |\phi_m|
 \end{aligned} \tag{4.3.5}$$

Thus by computing the bottom element of the solution vector to the tridiagonal system the residual norm can be computed using a simple multiplication, noting that  $\beta_{m+1}$  is computed at the  $m$ -th Lanczos step. In fact it is not even necessary to compute the solution to the tridiagonal system,  $f_m$  to obtain  $\phi_m$ . A recurrence for updating  $\phi_j$  from information at the previous step can be derived by considering a QR factorization of  $T_m$ . Consider first a factorized form of the tridiagonal matrix at step  $j - 1$

$$T_{j-1} = R_{j-1} U_{j-1} \tag{4.3.6}$$

where  $R_{j-1}$  is the product of all the rotation matrices up to and including step  $j - 1$  that transforms  $T_{j-1}$  to the upper triangular matrix  $U_{j-1}$ . Then the system of equations 4.2.3, for step  $j - 1$ , can now be written as

$$U_{j-1}f_{j-1} = \beta_1 R_{j-1}^T e_1^{(j-1)} \quad (4.3.7)$$

The last equation of the above can be written as

$$\tilde{\alpha}_{j-1} \phi_{j-1} = \gamma_{j-1} \quad (4.3.8)$$

Here tilde denotes all quantities that were modified by the all the previous rotations and  $\gamma_{j-1}$  is the bottom element that is created by these rotations. After an additional Lanczos step the last two equations of the new system is

$$\begin{bmatrix} \tilde{\alpha}_{j-1} & \tilde{\beta}_j \\ \beta_j & \alpha_j \end{bmatrix} \begin{bmatrix} \phi_{j-1} \\ \phi_j \end{bmatrix} = \begin{bmatrix} \gamma_{j-1} \\ 0 \end{bmatrix} \quad (4.3.9)$$

these equations are transformed again into an upper triangular form by premultiplication of equation 4.3.9 by the plane rotation matrix

$$\begin{bmatrix} c_j & -s_j \\ s_j & c_j \end{bmatrix}$$

where  $c_j = \cos \psi_j$  and  $s_j = \sin \psi_j$ . Here  $\psi_j$  is the rotation angle. The result of this matrix multiplication is

$$\begin{bmatrix} c_j \tilde{\alpha}_{j-1} - s_j \beta_j & c_j \tilde{\beta}_j - s_j \alpha_j \\ s_j \tilde{\alpha}_{j-1} + c_j \beta_j & s_j \tilde{\beta}_j + c_j \alpha_j \end{bmatrix} \begin{bmatrix} \phi_{j-1} \\ \phi_j \end{bmatrix} = \begin{bmatrix} c_j \gamma_{j-1} \\ s_j \gamma_{j-1} \end{bmatrix} \quad (4.3.10)$$

The value of the rotation angle is determined by setting the term in the lower triangular part equal to zero. Thus

$$\tan \psi_j = \frac{\beta_j}{\tilde{\alpha}_{j-1}} \quad (4.3.11)$$

and finally  $\phi_j$  can be obtained by performing the first step of the backsubstitution process.

$$\phi_j = \frac{s_j \gamma_{j-1}}{(s_j \tilde{\beta}_j + c_j \alpha_j)} \quad (4.3.12)$$

It should be pointed out that the rotation at this step effects the next off-diagonal term in the tridiagonal matrix and the effect is only on the super-diagonal term which explains the  $\tilde{\beta}_j$  term in equation 4.3.9. The final form of the updating algorithm can be arranged as follows

$$\begin{aligned} \tan \psi_j &= \frac{\beta_j}{\tilde{\alpha}_{j-1}} \\ \tilde{\alpha}_j &= \sin \psi_j \tilde{\beta}_j + \cos \psi_j \alpha_j \\ \tilde{\beta}_{j+1} &= \cos \psi_j \beta_{j+1} \\ \tilde{\gamma}_j &= \sin \psi_j \tilde{\gamma}_{j-1} \\ \phi_j &= \frac{\tilde{\gamma}_j}{\tilde{\alpha}_j} \end{aligned}$$

This is a very inexpensive step and therefore an attractive form of updating the bottom element of the solution vector. The residual norm,  $\rho_j$ , can then be computed from equation 4.3.5.

Contrary to the Conjugate Gradient method it is not necessary to update an approximate solution vector at each step, but only the residual norm has to be monitored. Once  $\rho_j$  falls below the specified tolerance the Lanczos algorithm is terminated and the small  $m$  degree tridiagonal system is solved. It is very important to note that if  $A$  is indefinite then  $T_j$  may also be indefinite and care must be taken in solving the small system. For this reason a stable equation solver such as the one used in the updating procedure for the residual norm should be used.

Finally the approximate solution  $x_m$  has to be computed. It is here that the Lanczos vectors are needed again to construct  $x_m$ . In the present implementation the Lanczos vec-

tors are put onto secondary storage as they are computed and at the end of the algorithm they are brought back one by one to form  $x_m$ .

#### 4.4. Derivation of Conjugate Gradients Method from Lanczos Algorithm

The conjugate gradients method may be developed directly from the Lanczos process. This will demonstrate the limitation of this method. The aim is to obtain a procedure for updating of the solution vector  $x_j$  without retaining all the Lanczos vectors.

Consider first the case where  $A$  is positive definite, then the resulting tridiagonal matrix will also be positive definite and hence its cholesky factorization exists.

$$T_j = L_j L_j^T \quad (4.4.1)$$

Here  $L_j$  is a lower bidiagonal matrix with positive diagonal terms.

$$L_j = \begin{bmatrix} \delta_1 & & & \\ \mu_2 & \delta_2 & & \\ & \mu_3 & \ddots & \\ & & \ddots & \delta_{j-1} \\ & & & \mu_j & \delta_j \end{bmatrix} \quad (4.4.2)$$

The main obstacle to updating  $x_j$  is the fact that the solution vector  $f_j$  to equation 4.2.3 changes fully after each Lanczos step and therefor  $Q_j f_j$  can not be simply accumulated recursively. This difficulty can be overcome by defining

$$g_j = L_j^T f_j \quad (4.4.3)$$

and

$$P_j = Q_j L_j^{-T} \quad (4.4.4)$$

The vectors in  $P_j$  are the conjugate vectors. Then the approximation  $x_j$  becomes

$$x_j = P_j g_j \quad (4.4.6)$$

Further, the weak form equation now becomes

$$L_j g_j = e_1^{(j)} \beta_1 \quad (4.4.7)$$

It is apparent from the bidiagonal structure of  $L_j$ , equation 4.4.2, the  $g_j$  can be updated very easily since the solution of equation 4.4.7 is obtained by a simple forward reduction. Therefor  $g_{j+1}$  can be written in terms of  $g_j$  as

$$g_{j+1} = \begin{pmatrix} g_j \\ \eta_{j+1} \end{pmatrix} \quad (4.4.8)$$

where

$$\eta_{j+1} = -\mu_{j+1} g_j / \delta_{j+1} \quad (4.4.9)$$

Thus the approximation  $x_{j+1}$  can be obtained from  $x_j$  using equation 4.4.6 in arranged in the following form

$$x_{j+1} = x_j + \eta_{j+1} p_{j+1} \quad (4.4.10)$$

and finally the  $p$ -vectors can be obtained from their definition. Thus

$$p_{j+1} = \frac{1}{\delta_{j+1}} (q_{j+1} - \mu_{j+1} p_j) \quad (4.4.11)$$

This is equivalent to the method of conjugate gradients. The approach here is computationally a little different, but the role of the cholesky decomposition becomes apparent, with the necessary requirement that  $A$  be positive definite to ensure numerical stability.

If  $A$  is an indefinite symmetric matrix, the algorithm may still be carried out, with some success, but due to instability it can no longer be numerically reliable.

#### 4.5. Selective Orthogonalization

A detailed account of the behaviour of the Lanczos algorithm in the presence of roundoff and of selective orthogonalization is available in Parlett [18]. Therefore in this section only the basic facts about selective orthogonalization will be presented.

Let  $\epsilon$  be the smallest number in the computer such that  $1 + \epsilon > 1$ . This is known as the unit roundoff error. The basic equations 4.1.5.a and 4.1.5.c are now perturbed by roundoff. Although for each  $j$  4.1.5.b is only slightly perturbed, the relations 4.1.5.a and 4.1.5.c completely fail after a certain number of steps depending on  $\epsilon$  and on  $A$ . The Lanczos vectors, which are orthogonal in exact arithmetic, not only lose their orthogonality, but even become linearly dependent.

For a long time, as a remedy against this loss of orthogonality, it was suggested to reorthogonalize each new  $q_j$  against all previous Lanczos vectors, which is of course very expensive. The results of Paige [15] gave some better insight in the way orthogonality is lost, and provide the theoretical basis for selective orthogonalization.

In order to explain Paige's results it is necessary to introduce certain quantities which are not of direct interest when solving a linear system of equations. Let  $\theta_i^{(j)}$  be the eigenvalues and  $s_i^{(j)}$  be the corresponding normalized eigenvectors of  $T_j$

$$T_j s_i^{(j)} = s_i^{(j)} \theta_i^{(j)} \quad i = 1, \dots, j \quad (4.5.1)$$

From these one can compute the Ritz vectors  $y_i^{(j)}$  by

$$y_i^{(j)} = Q_j s_i^{(j)} \quad (4.5.2)$$

These quantities change at each Lanczos step. In subsequent discussions where there is no confusion possible the superscripts will be dropped. The pairs  $(\theta_i^{(j)}, y_i^{(j)})$   $i=1, \dots, j$  are approximate eigenpairs of  $A$ . The quality of this approximation can be determined by considering the residual  $\|Ay_i - y_i \theta_i\|$ . Applying (4.5.2) one finds that, to within roundoff,

$$\|Ay_i - y_i \theta_i\| \leq \beta_j + \|E_j\| \quad (4.5.3)$$

where  $E_j$  is an error matrix and  $\beta_j \equiv \beta_{j+1} |s_j|$  and  $s_j \equiv e_j^T s_j$ , i.e., the bottom element of the corresponding eigenvector of  $T_j$ . The  $\beta_j$  play an important role in the process of selective orthogonalization.

Now everything is ready to state one of the most important consequences of Paige's work, which can be summarized as follows:

- Loss of orthogonality among the Lanczos vectors is equivalent to the convergence of a Ritz pair.

In other words, if one of the  $\beta_{ji}$  becomes small, the corresponding Ritz pair converges to an eigenpair of  $A$  and the Lanczos vector  $q_{j+1}$  loses its orthogonality to  $q_1, \dots, q_{j-2}$ . But more is known about the way  $q_{j+1}$  behaves, it is tilted towards  $y_i^{(j)}$  while it is retaining its previous level of orthogonality to all the other Ritz vectors,  $y_k^{(j)}, k \neq i$ . In order to maintain a certain level of orthogonality it is therefore only necessary to orthogonalize the new  $q_{j+1}$ , or equivalently the unnormalized  $r_j$  against  $y_i^{(j)}$  when the  $\beta_{ji}$  become smaller than some threshold. So first compute  $r'_j$  (compare 4.1.3) by

$$r'_j \equiv Aq_j - q_j\alpha_j - q_{j-1}\beta_j \quad (4.5.4)$$

as usual, then check if any  $\beta_{ji}$  is small. If so then compute the corresponding  $y_i^{(j)}$  and orthogonalize  $r'_j$  against it obtaining

$$r_j \equiv r'_j - y_i^{(j)}\xi_i^{(j)} \quad (4.5.5)$$

where

$$\xi_i^{(j)} \equiv \frac{y_i^{(j)T} r'_j}{\|y_i^{(j)}\|^2} \quad (4.5.6)$$

This requires the formation of  $y_i = Q_j s_i$ . The payoff for the expensive calculation comes later when subsequent Lanczos vectors, say  $q_{j+15}$  and  $q_{j+30}$  need to be orthogonalized against  $y_i$ , since a certain number  $p$  of Ritz vectors are also put in secondary storage and need not be formed again. The question of how many Ritz vectors to keep, and whether to keep them in core in case the optimal number is small is still under investigation.

Finally it should be mentioned that the threshold for which a  $\beta_{ji}$  is considered to be small is set to be  $\sqrt{\epsilon}\|T_j\|$  in the present implementation. This choice was based on the

computational experience for the eigenvalue problem [20]. An analysis in [20] shows that this guarantees a certain level of orthogonality among the Lanczos vectors.

To demonstrate the loss of orthogonality, in Figure 4.2, the variation in the residual norm is plotted for the simple Lanczos process when no action is taken to maintain orthogonality and for the case when selective orthogonalization is performed to maintain a level of orthogonality among the Lanczos vectors. The conjugate gradient method and other methods based on the Krylov vectors also suffer this effect. It should be noted that the loss of orthogonality does not stop convergence, it only delays it.

#### 4.6. Lanczos with Selective Orthogonalization

The final form of the Lanczos algorithm which maintains a certain level of orthogonality is summarized in the following table.

1. Loop: for  $j = 1, 2, \dots$ 
  - 1.1 Take a simple Lanczos step
  - 1.2 Update residual norm  $\rho_j$
  - 1.3 Check: if  $\rho_j \leq \text{tolerance}$  end the loop
  - 1.4 Update the eigenvalues of  $T_j$
  - 1.5 Compute  $\beta_{ji} = \beta_{j+1} s_{ji}$
  - 1.6 Check: if  $\beta_{ji} \leq \sqrt{\epsilon} \|T_j\|$   
then compute the  $i$ -th eigen vector and orthogonalize
2. Solve  $T_j f_j = e_j \beta_1$
3. Assemble solution:  $x_j = Q_j f_j$

The above formulation requires a total of 5  $n$ -vectors and 10  $j$ -vectors workspace. This includes two  $n$ -vectors containing the right hand side and the initial guess, two  $n$ -vectors for  $r_j$  and  $q_j$  and a  $n$ -vector for use in construction of Ritz vectors.



#### 4.7. Selective Reorthogonalization

The procedure described in the previous section is not the only way of maintaining orthogonality among the Lanczos vectors. Recently H. Simon [22] has proposed an alternative technique based on a simple recursion formula. This recursion is derived as follows.

Define

$$\omega_{i,j} = \mathbf{q}_i^T \mathbf{q}_j \quad (4.7.1)$$

and multiplying equation 4.1.3 by  $\mathbf{q}_i^T$  and using the above definition we get

$$\mathbf{q}_i^T \mathbf{A} \mathbf{q}_j = \alpha_j \omega_{i,j} + \beta_j \omega_{i,j-1} + \beta_{j+1} \omega_{i,j+1} \quad (4.7.1.a)$$

a similar equation can be obtained when the same steps are applied to the  $i$ -th Lanczos step.

$$\mathbf{q}_j^T \mathbf{A} \mathbf{q}_i = \alpha_i \omega_{j,i} + \beta_i \omega_{j,i-1} + \beta_{i+1} \omega_{j,i+1} \quad (4.7.1.b)$$

subtracting the above two relations to eliminate the terms involving  $\mathbf{A}$  results in the recursion

$$\beta_{j+1} \omega_{j+1,i} = \beta_{i+1} \omega_{j,i+1} + (\alpha_i - \alpha_j) \omega_{j,i} + \beta_i \omega_{j,i-1} - \beta_j \omega_{j,i-1} \quad (4.7.3)$$

This recursion holds for  $j = 2, 3, \dots$  and  $1 \leq i \leq j-1$ . To start the three step recursion it is necessary to assume that

$$\omega_{j,j} = 1 \quad j = 1, 2, \dots \quad (4.7.4.a)$$

and

$$\omega_{j,j-1} = \epsilon \quad j = 2, 3, \dots \quad (4.7.4.b)$$

Using the above recursion the loss of orthogonality can be monitored during the Lanczos process simply by looking at the magnitude of  $\omega_{j,i}$ . Whenever  $|\omega_{j,i}|$  is larger than the threshold value,  $\sqrt{\epsilon} \|\mathbf{T}_j\|$ , orthogonality is considered to be lost against  $\mathbf{q}_i$  and reorthogonalization must be performed. Lanczos algorithm with selective reorthogonalization then becomes:

1. Loop: for  $j = 1, 2, \dots$ 
  - 1.1 Take a simple Lanczos step
  - 1.2 Update residual norm  $\rho_j$
  - 1.3 Check: if  $\rho_j \leq \text{tolerance}$  end the loop
  - 1.4 Update  $\omega_{i,j}$
  - 1.5 Check: if  $\omega_{i,j} \geq \sqrt{\epsilon} \|T_j\|$

then orthogonalize  $q_{j+1}$  against previous Lanczos vectors.

2. Solve  $T_j f_j = e_j \beta_1$

3. Assemble solution:  $x_j = Q_j f_j$

This algorithm has a number of desirable advantages over the previous approach.

- (i) The implementation of this recursion is very easy and the resulting code is considerably simpler.
- (ii) This approach requires less storage space.
- (iii) There is no need for computing the eigenvalues and eigenvectors of the tridiagonal system.

Initial experiments with the two approaches indicate that both methods flag the loss of orthogonality at the same step. Therefore due to the above advantages the selective reorthogonalization technique is a better way of maintaining orthogonality for purpose of equation solving.

#### 4.8. Starting Vector

The choice of a good starting vector can reduce the number of Lanczos steps the algorithm takes to converge to the solution. In fact if the true solution is on hand the Lanczos algorithm terminates after one step. However it is not often that one has the exact solution or even a good approximation. Nevertheless the question of starting vectors still remains. Consider the extreme case where the right hand side is an eigenvector. Then any starting

vector other than zero or the right hand side will take more than one step to converge, whereas, a zero starting vector will take exactly one step to obtain the true solution. This can be seen by performing the algorithm on the equation  $Ax = z_i$ , where  $z_i$  is the  $i$ -th eigenvector with corresponding eigenvalue  $\lambda_i$ . Then

Setting  $q_0 = 0$

$r_0 = z_i$  and  $\beta_1 = 1$

$$(1) \ q_1 \leftarrow r_0/\beta_1 = z_i$$

$$(2) \ u_1 \leftarrow Aq_1 = \lambda_i z_i$$

$$(3) \ r_1 \leftarrow u_1 - \beta_1 q_0 = \lambda_i z_i$$

$$(4) \ \alpha_1 \leftarrow q_1^T r_1 = \lambda_i$$

$$(5) \ r_1 \leftarrow r_1 - \alpha_1 q_1 = 0$$

$$(6) \ \beta_2 \leftarrow 0$$

and the residual norm will be zero.

In general the right hand side has components along a certain number of eigenvectors and introduction of any starting vector that has components along any additional eigenvectors will delay convergence. Therefore, it can be said that unless special knowledge about the solution is available, the best starting vector is the zero vector. The influence of a random starting vector on the residual norm during a Lanczos run can be seen in Figure 4.3 for the simple beam problem.

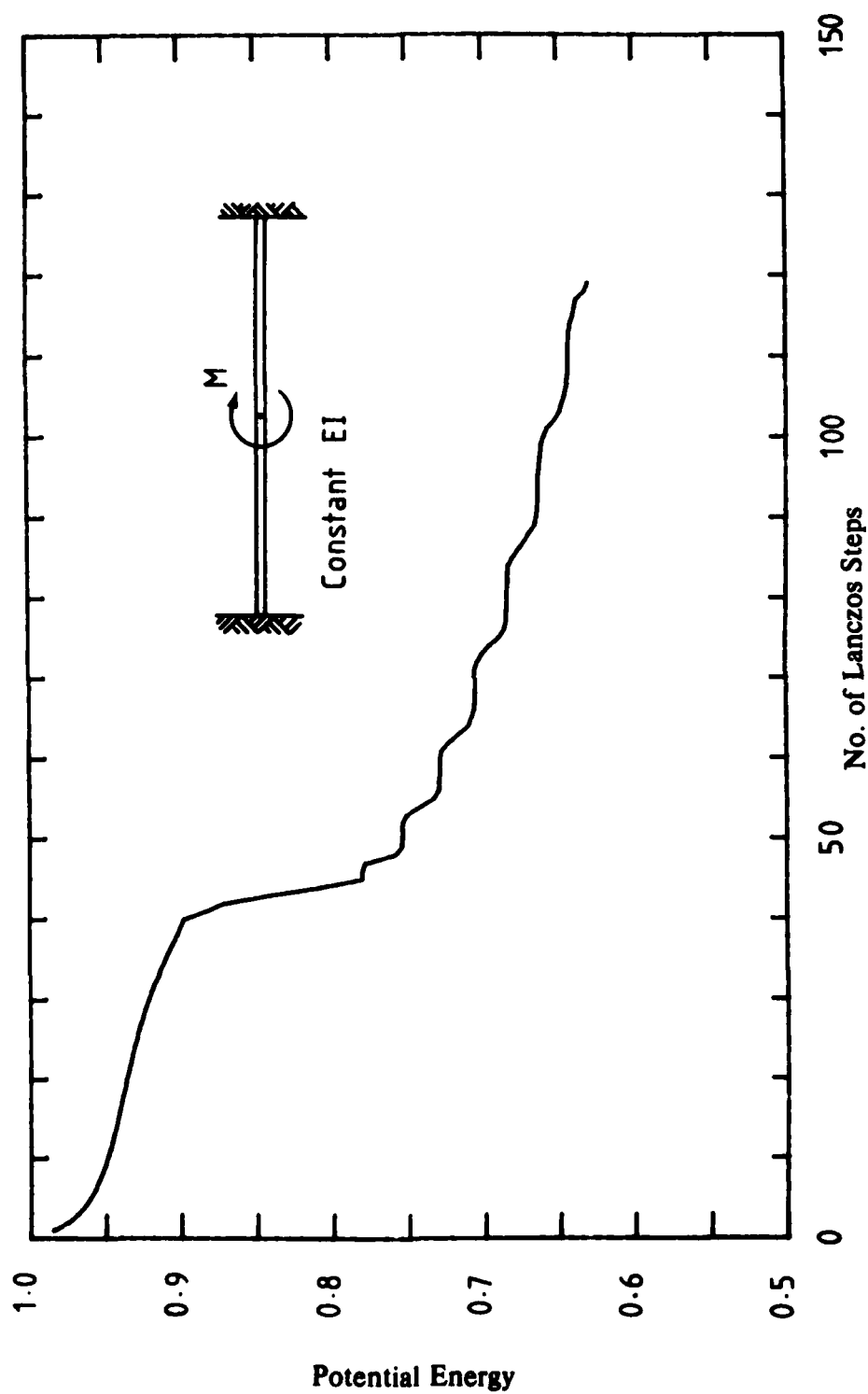


Figure 4.1 Monotonic Reduction of the Potential Energy of the System in the Lanczos Algorithm.

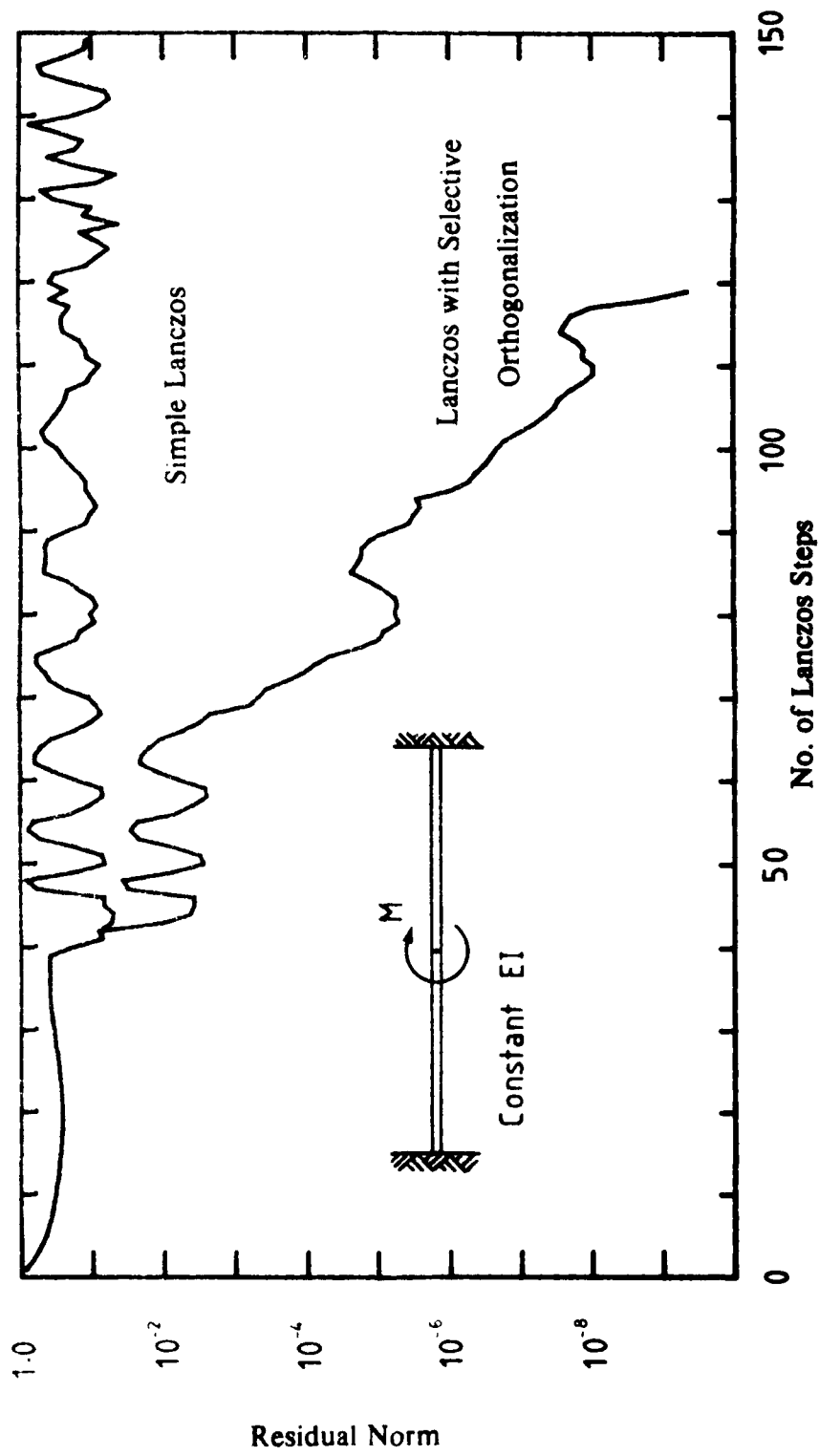


Figure 4.2 Effect of Selective Orthogonalization on the Reduction of the Residual Norm in the Lanczos Algorithm.

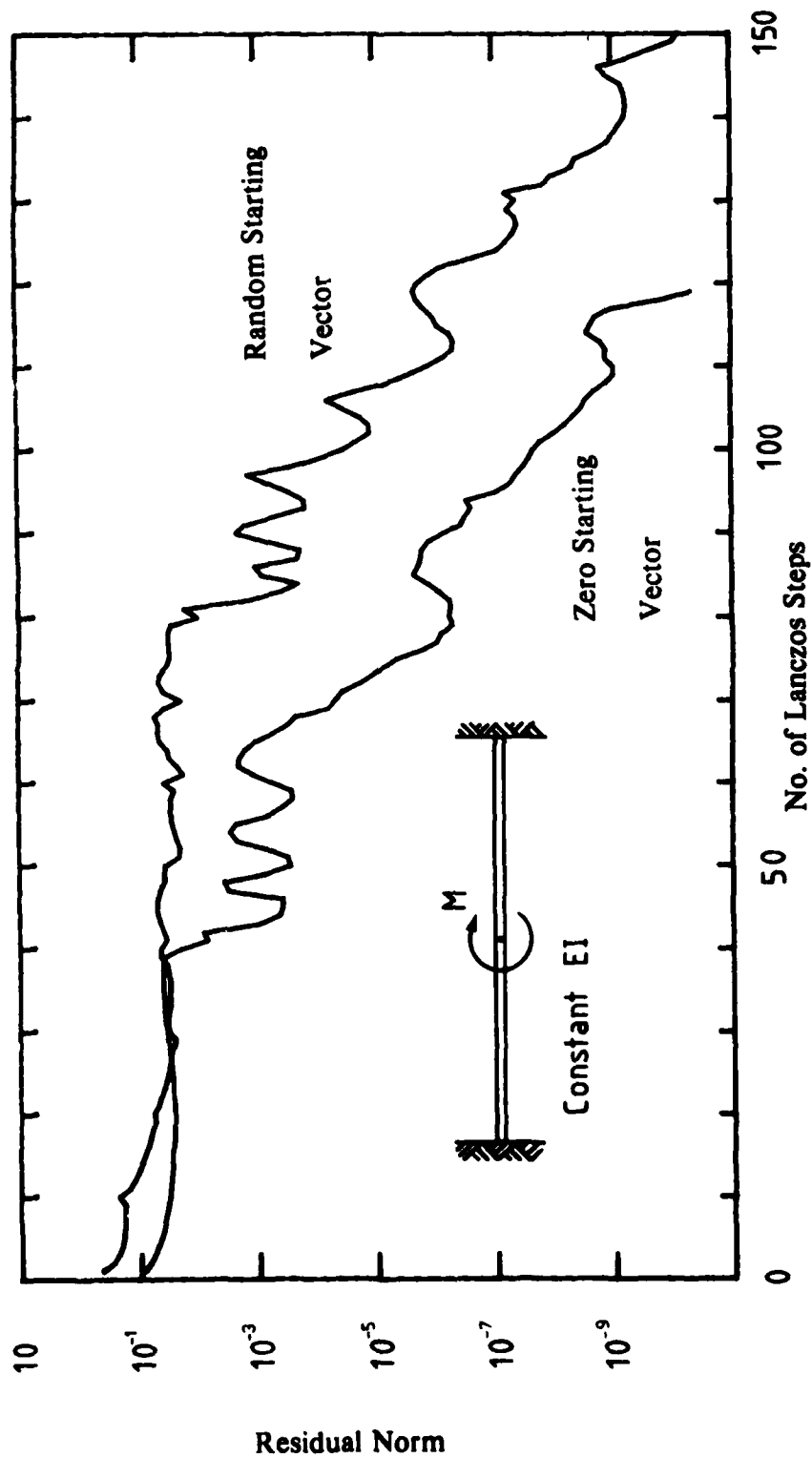


Figure 4.3 Effect of Different Starting Vectors on the Reduction of the Residual Norm in the Lanczos Algorithm.

## 5. NEWTON-LANCZOS METHOD

Generally in nonlinear structural analysis, the solution to the nonlinear equations is required at a number of load levels to generate a complete equilibrium history. The BFGS method described in section 3.3 attempts to make use of the information at a previous iterate, by updating the factorization of the tangent matrix, to find the next solution. This updating procedure imposes a limitation on the solution algorithm. That is:

The BFGS can not be applied to nonlinear problems where the tangent matrix becomes indefinite during the iteration process.

The object of the BFGS algorithm was to reduce the overall number of matrix factorizations but this could only be achieved at the cost of limiting the class of problems for which the method can be useful.

In addition to the short comings of Newton type methods described in section 3, all these methods including Newton's method have the following limitations.

- (1) These methods are not globally convergent, i.e not all starting vectors converge to the solution.
- (2) The algorithm can not be applied to the cases when the tangent matrix becomes singular at some  $x_k$ .

Here global convergence means that if  $f(x) = 0$  has a unique solution then all starting vectors converge to this solution.

The object of the proposed algorithm is

- (i) Remove the above two restrictions.
- (ii) Reduce the number of matrix factorization.

### 5.1. Description of the Algorithm

The Newton-Lanczos algorithm is based on the simple idea that so long as the approximate solution is far from the exact solution, it is not necessary to solve the linearized

problem exactly, but only a modest level of accuracy is sufficient.

The algorithm consists of choosing an initial guess  $x_0$  and repeating the following steps for  $k = 0, 1, 2, \dots$ , until convergence is achieved:

- (1) given  $x_k$  compute  $f(x_k)$
- (2) compute the tangent matrix  $A_k$
- (3) find some  $d_k$  such that  $\frac{\|A_k d_k + f(x_k)\|}{\|f(x_k)\|} \leq \eta_k$
- (4) compute  $s_k$  from a line search
- (5) update  $x_{k+1} = x_k + s_k d_k$
- (6) test for convergence and terminate if converged

Here  $\eta_k$  is some parameter such that  $\eta_k < 1$ . The only difference between this algorithm and Newton's method is in step 3. At this step the Lanczos algorithm as described in section 4 is used to obtain  $d_k$  and  $\eta_k$  is the specified solution tolerance for the Lanczos algorithm.

It is evident from the description of the algorithm that if  $A_k$  is singular at some step  $k$  then the process need not be terminated as long as a  $d_k$  can be found that can achieve the required tolerance. Under this circumstance it is desirable to have small components of the eigenvectors corresponding to the zero eigenvalues along  $d_k$ . The behavior of the Lanczos algorithm is such that the process takes close to  $n$  steps, where  $n$  is the number of unknowns, before it introduces any component of these eigenvectors into the solution vector. The approximate solution is obtained from a linear combinations of the Krylov vectors and it is the matrix-vector multiplications that annihilates components of the eigenvectors corresponding to zero eigenvalue. Therefore, it is possible to obtain an approximate solution without introducing large components of the undesirable eigenvectors.

The cost of step 3 can be reduced by obtaining a good initial guess to the solution of the new system. This can be done by solving the weak form of the linearized equation at a



previous step  $i$ ,

$$\mathbf{T}_j^{(i)} \mathbf{g}_j = \mathbf{h}_j$$

with the appropriate right hand side. The projection of the residual vector,  $\mathbf{f}(\mathbf{x}_k)$ , on the Lanczos vectors of the previous step is identically this right hand side, That is

$$\mathbf{h}_j = \mathbf{Q}_j^{(i)T} \mathbf{f}(\mathbf{x}_k)$$

## 5.2. Convergence of Newton-Lanczos

Newton-Lanczos algorithm has very attractive convergence properties. The solution part of the algorithm can be rearranged to get

$$\|\mathbf{A}_k \mathbf{d}_k + \mathbf{f}(\mathbf{x}_k)\| \leq \eta_k \|\mathbf{f}(\mathbf{x}_k)\| \quad (5.2.1)$$

The convergence of Newton-Lanczos method can be proved provided the following assumptions hold.

- (a) The solution  $\mathbf{x}^*$  to  $\mathbf{f}(\mathbf{x})$  is unique.
- (b) The Jacobian matrix exists.
- (c)  $\mathbf{f}'(\mathbf{x}^*)$  is nonsingular.

The Taylor expansion about  $\mathbf{x}_k$  can be written as

$$\mathbf{f}(\mathbf{x}_{k+1}) = \mathbf{f}(\mathbf{x}_k) + \mathbf{A}_k \mathbf{d}_k + \mathbf{e}(\mathbf{d}_k) \quad (5.2.2)$$

where  $\mathbf{e}(\mathbf{d}_k)$  includes all the higher order terms in  $\mathbf{d}_k$ . Hence

$$\|\mathbf{f}(\mathbf{x}_{k+1})\| \leq \|\mathbf{f}(\mathbf{x}_k) + \mathbf{A}_k \mathbf{d}_k\| + \|\mathbf{e}(\mathbf{d}_k)\| \quad (5.2.3)$$

Using equation 5.2.1 this can be modified to

$$\|\mathbf{f}(\mathbf{x}_{k+1})\| \leq \eta_k \|\mathbf{f}(\mathbf{x}_k)\| + \|\mathbf{e}(\mathbf{d}_k)\| \quad (5.2.4)$$

At this stage there are two possible cases that need be considered:

- (1)  $\eta_k \|\mathbf{f}(\mathbf{x}_k)\| \leq \|\mathbf{e}(\mathbf{d}_k)\|$

This implies that  $\|\mathbf{f}(\mathbf{x}_{k+1})\| < 2\|\mathbf{e}(\mathbf{d}_k)\| = O(\|\mathbf{d}_k\|^2)$  which in turn implies at least

superlinear convergence.

$$(2) \|e(d_k)\| < \eta_k \|f(x_k)\|$$

This inequality can be stated as an equality of the form

$$\|e(d_k)\| = t\eta_k \|f(x_k)\| \quad (5.2.5)$$

where  $t$  is a positive scalar such that  $t < 1$ , then equation 5.2.4 can be written as

$$\|f(x_{k+1})\| < \gamma \|f(x_k)\| \quad (5.2.6)$$

where  $\gamma = (t+1)\eta_{\max}$  is a scalar. Choosing  $\eta_{\max}$  such that  $\eta_k < \eta_{\max} < \gamma < 1$  for all values of  $k$  Then

$$\lim_{k \rightarrow \infty} f(x_k) = 0 \quad (5.2.7)$$

These results show that the algorithm is at least linearly convergent, i.e

$$\|x_{k+1} - x^*\| \leq \alpha \|x_k - x^*\| \quad \alpha < 1 \quad (5.2.8)$$

It can also be observed directly from the algorithm that by setting  $\eta_k = \epsilon$ , the computer precision, the algorithm will achieve a quadratic rate of convergence.

In general the rate of convergence is superlinear.

$$\|x_{k+1} - x^*\| \leq \alpha \|x_k - x^*\|^{q(\eta_k)} \quad (5.2.9)$$

where  $1 \leq q(\eta_k) \leq 2$  and  $\lim_{\eta \rightarrow 0} q(\eta) = 2$ . The superlinear rate of convergence can be guaranteed by choosing  $\eta_k$  such that  $\eta_k \rightarrow 0$  as  $x_k \rightarrow x^*$ . This can be achieved very simply by making  $\eta_k$  dependent on  $\frac{\|f(x_k)\|}{\|f(x_0)\|}$ . In the present implementation this dependence was chosen as

$$\eta_k = \eta_0 \left[ \frac{\|f(x_k)\|}{\|f(x_0)\|} \right]^{1.5}$$

based on some computational experiments.

### 5.3. Preconditioning of The Linear Equations

The use of preconditioning of the linear system of equations can often speed up convergence of the Lanczos algorithm considerably. The number of Lanczos steps required to reduce the residual norm by a factor  $\eta$  depends, in general, on the distribution of the eigenvalues of  $A$ . An upper bound,  $p$ , to this number can be obtained from following relation

$$\eta = \left[ \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right]^p \quad (5.3.1)$$

where  $\kappa$  is the condition number of  $A$ ,  $\kappa = \lambda_n/\lambda_1$ . This upper bound is too crude to have any practical use for estimating the number of Lanczos steps but it can be seen that when the condition number is unity Lanczos algorithm requires only one step to obtain the exact solution. The equation 4.1.1 can now be modified to

$$C^{-1}AC^{-T}\hat{x} = \hat{r} \quad (5.3.2)$$

where

$$\hat{x} = C^T x^c \quad (5.3.3)$$

and

$$\hat{r} = C^{-1}r_0 \quad (5.3.4)$$

where  $C$  is any matrix with an inverse that can be computed easily such as a triangular matrix. The object now is to solve 5.3.2. The number of Lanczos steps can be reduced substantially by choosing  $C$  in such a way as to reduce the condition number of the product matrix of equation 5.3.2.

The present implementation of the Newton-Lanczos algorithm uses the Cholesky factor of the tangent matrix at some previous step for  $C$ . In this way the information is passed on from one iterate to the next and therefore reduces the number of factorizations. However, factorizations can be eliminated completely if a good approximation to the preconditioning matrix  $C$  can be obtained. Moreover, the preconditioning can be updated from one

step to the next using an updating procedure similar to the one used for the BFGS algorithm.

#### 5.4. Numerical Results

##### *Example I*

This example, figure 5.1, is a simple truss structure with a peculiar property. The initial tangent matrix is singular. The purpose of this example is to demonstrate the global convergence of the Newton-Lanczos algorithm. A linearly elastic, large displacement truss element was used to model the structure which has a total of 8 degrees of freedom.

Using a zero starting solution, all the methods described in section 3, including Newton's method, diverged immediately after the initial iteration when applied to this problem. This is due to the fact that a direct equation solver is used to solve the linearized problem and since the tangent matrix is singular the solution vector will be an eigenvector corresponding to a zero eigenvalue with a large magnitude. This is like performing one step of inverse iteration method with a singular matrix which always converge to the eigenvector with smallest eigenvalue. No line search was used for Newton's method or Modified Newton method, however the BFGS algorithm which has line search also failed to converge.

Newton-Lanczos algorithm, with no preconditioning, when applied to this problem converged to the correct solution. However the convergence was only linear because of the singularity. The convergence rate improved as the approximate solution moved away from the singularity. Despite the low rate of convergence, the fact that the correct solution was obtained shows that the domain of attraction of Newton-Lanczos is larger than that of Newton's method. The load-displacement curve of figure 5.2 gives an indication of the stiffening behavior of the structure.

##### *Example II*

For a second example, we consider the shear loading of a block of rubber, as shown in Figure 5.3. The nonlinear elastic behavior of the rubber was modeled using a simplification of the general Rivlin model, originally proposed by Mooney [4]. The nonlinearities are present throughout the loading process. A nine node plane stress Lagrangian element, with penalty formulation [21], was used to discretize the problem to a total of 60 degrees of freedom.

The four basic methods were tested on this problem and the results are tabulated in Table 5.1. The convergence of the BFGS quasi-Newton method was disappointing since it required 11 iterations to converge at each step while the modified Newton which has only a linear convergence rate, required 14 iteration to converge.

The total number of factorizations required by each method is listed below.

Newton	9
BFGS	3
Modified Newton	3
Newton-Lanczos	1

The control of the rate of convergence for Newton-Lnaczos can be seen, Figure 5.4, when the method was used, without any preconditioning, to solve this problem with different  $\eta_0$  values.

### *Example III*

As a final example a parabolic beam illustrated in Figure 5.5 was modeled using 4 node, plane stress, quadrilateral elements. The nonlinearities enter the problem when the beam comes into contact with a rigid boundary. The contact condition is enforced by Lagrange multiplier constraints on those nodes indicating penetration and/or compressive contact force. The contact problem is of the same type as the problem of constrained optimization.

Both The BFGS and the Modified Newton method failed to converge when applied to this problem. This is because the tangent matrix becomes indefinite during the iteration process and the BFGS updates are incapable of accounting for the indefiniteness introduced by the contact elements.

The structure has stiffening characteristic during the loading process and a softening characteristic during the unloading phase. The total number of degrees of freedom is 870. The cost comparisons for Newton's method and Newton-Lanczos method are tabulated in Table 5.2. The total number of factorization required by each method is as follows:

Newton	13
Newton-Lanczos	2

### 5.5. Closure

The motivation for the development of *quasi-Newton methods* was reduction in the overall cost of the analysis, using fewer factorization steps. However, the use of quasi-Newton methods to solve nonlinear finite element equations has been restricted to problems with positive definite Jacobians.

The Newton-Lanczos algorithm described herein not only requires less factorization steps than either the quasi-Newton method or the modified Newton method but also has the following advantages over its rivals:

- (1) The method is globally convergent.
- (2) The rate of convergence can be controlled through a specified tolerance.
- (3) The algorithm does not terminate in the case of a near singular Jacobian matrix.
- (4) The algorithm can be applied to indefinite systems.

In the present implementation of Newton-Lanczos, the Jacobian matrix is decomposed into its Cholesky factors in order to get a preconditioning matrix for the Lanczos algorithm. This factorization is avoided if an effective preconditioning matrix can be

constructed by other means, resulting in an additional reduction in the overall storage requirement of the algorithm. Further, the preconditioning matrix can then be updated using a sparse updating procedure similar to the one used by BFGS algorithm, for use in the subsequent iterations.

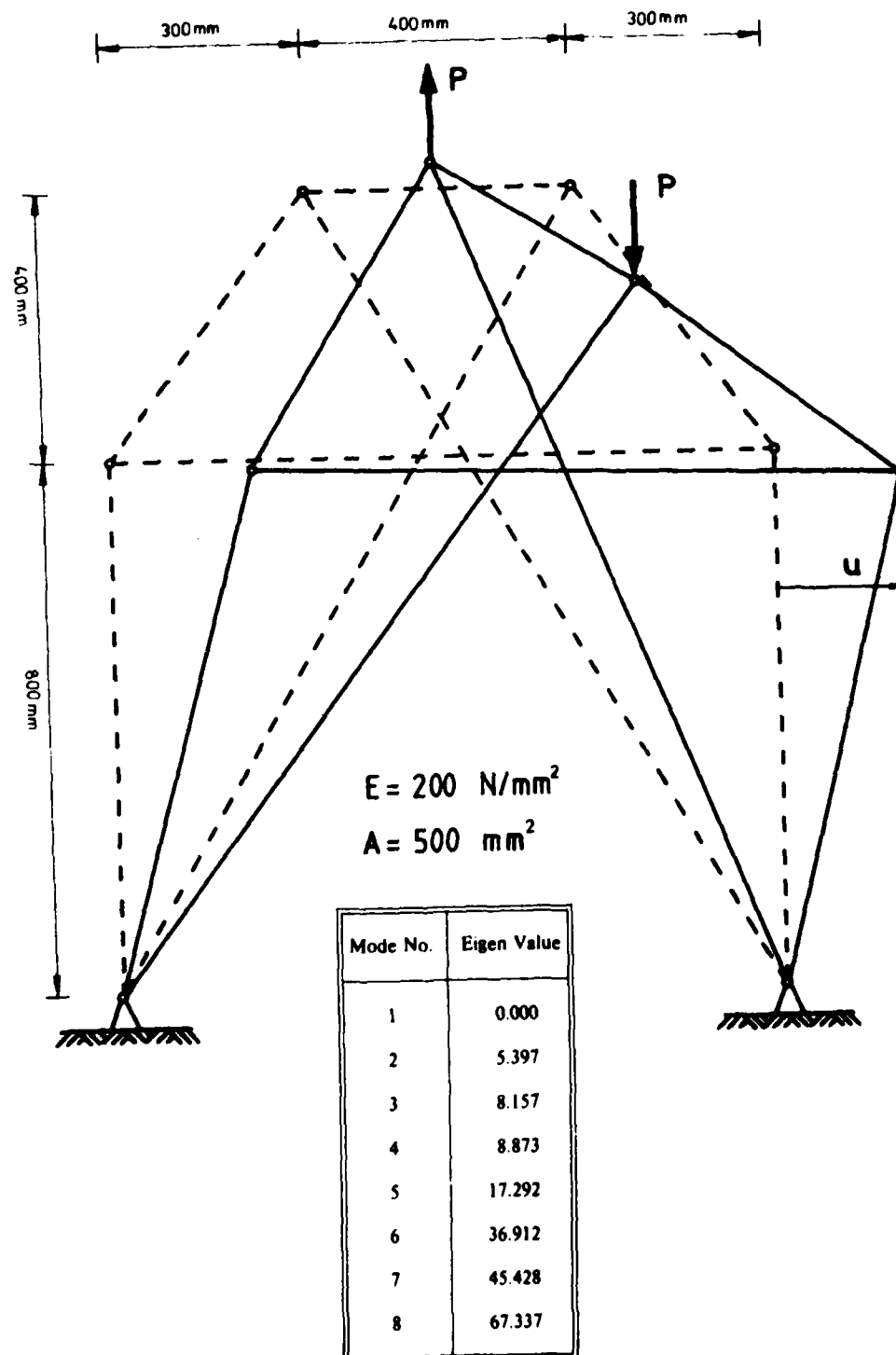


Figure 5.1 Properties of the Truss Structure used in Example I.



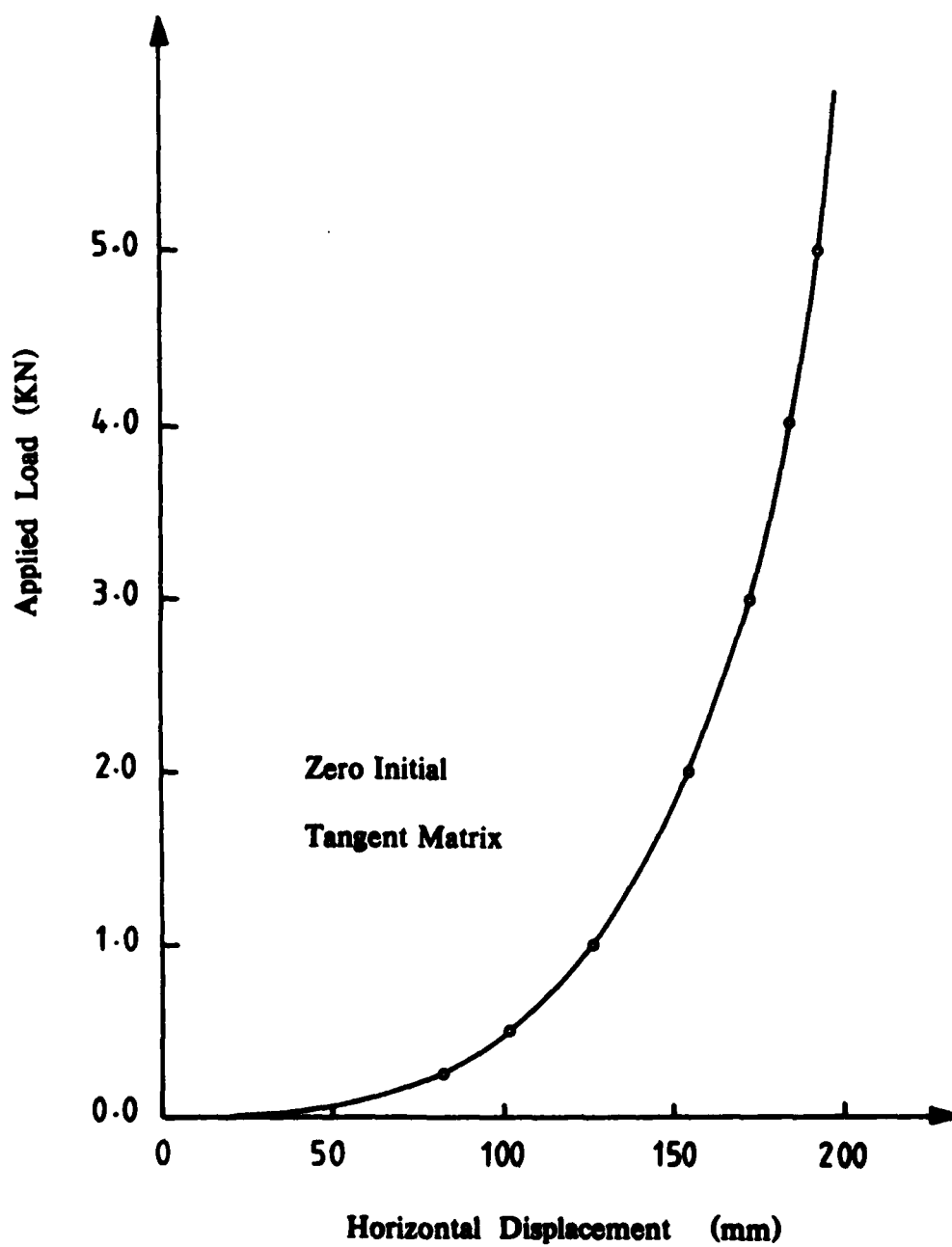
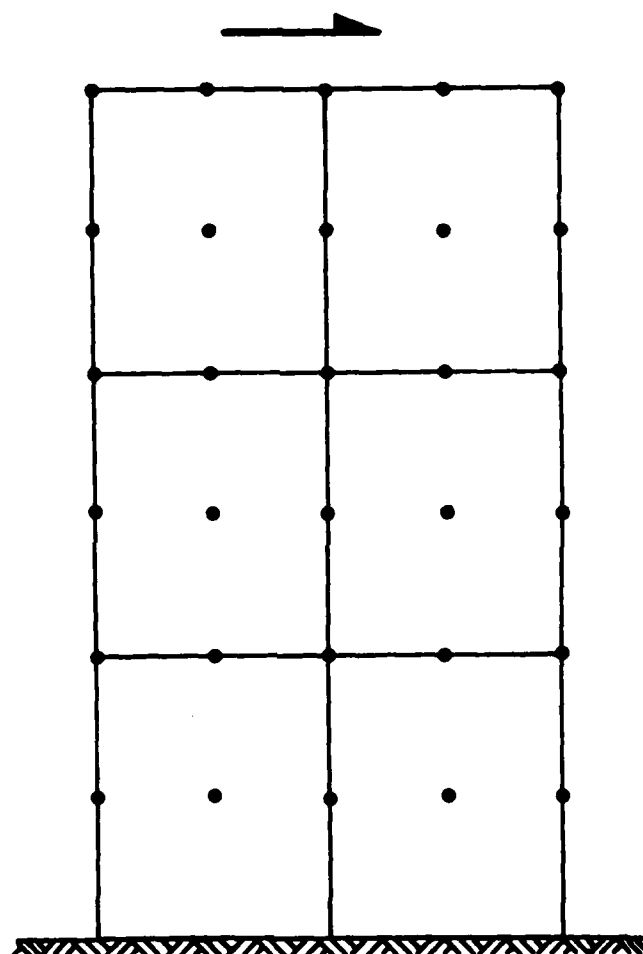


Figure 5.2 Load-Displacement Characteristic of Example I.

Nonlinear Elastic Material



Nine Node Lagrangian Elements

Figure 5.3 Problem of a Rubber Block under Shear used in Example II.

Load level	Method	No. of Iteration	No. of LU Factor.	No. of Function Evaluations	No. of Matrix-Vector ops.
1	Newton	3	3	4	3
	Mod. Newton	14	1	15	14
	BFGS	11	1	12	11
	N-L $\eta_0 = 10^{-3}$	4	1	5	25
	N-L $\eta_0 = 10^{-5}$	3	1	4	15
2	Newton	3	3	4	3
	Mod. Newton	14	1	15	14
	BFGS	11	1	12	11
	N-L $\eta_0 = 10^{-3}$	3	0	4	18
	N-L $\eta_0 = 10^{-5}$	3	0	4	24
3	Newton	3	3	4	3
	Mod. Newton	14	1	15	14
	BFGS	11	1	12	11
	N-L $\eta_0 = 10^{-3}$	3	0	4	18
	N-L $\eta_0 = 10^{-5}$	3	0	4	32

Table 5.1 Cost comparisons of different methods using

Example II ( Rubber Block ).

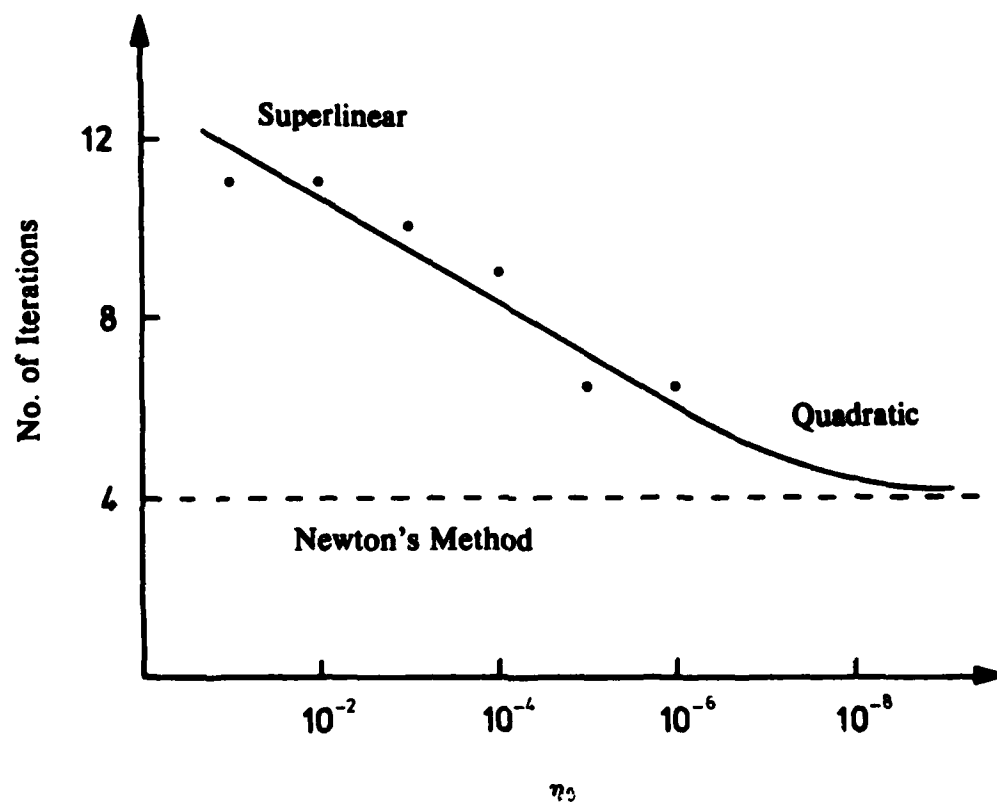


Figure 5.4 Effect of  $\eta_0$  on the Rate of Convergence of the Newton-Lanczos algorithm.

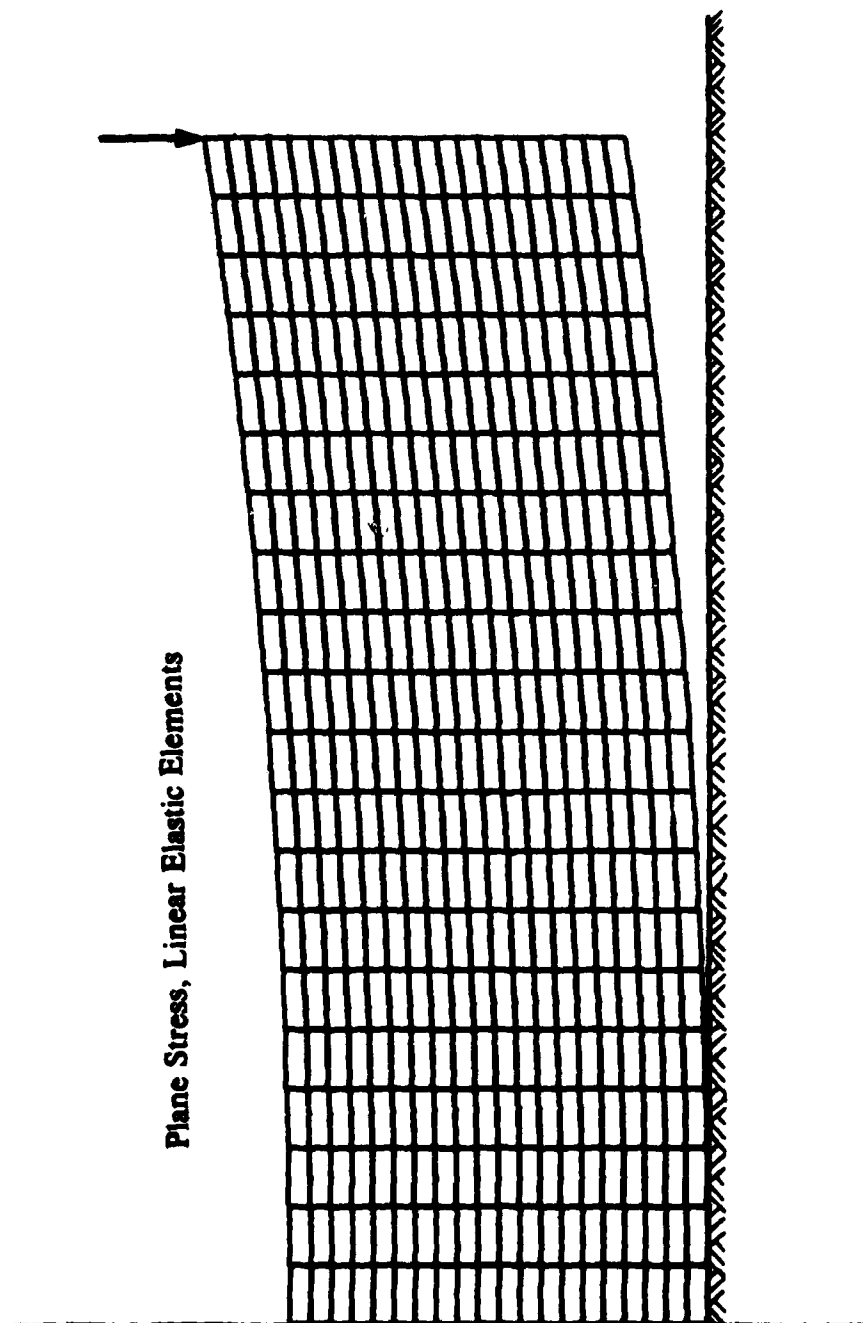


Figure 5.5 Contact Problem of a Parabolic Beam used in Example III.

Load (KN)	No. of Iterat.	Newton-Lanczos Method $\eta_0 = 10^{-7}$			Newton's Method	
		No. of LU Factor.	No. of Lanczos steps	C.P.U. time (sec.)	No. of LU Factor.	C.P.U. time (sec.)
0.4	3	2	6	64.01	3	67.22
0.6	2	0	3	18.11	2	44.43
0.9	2	0	7	28.29	2	43.69
0.6	2	0	8	33.18	2	44.45
0.4	2	0	4	22.16	2	44.45
0.0	2	0	6	31.09	2	44.32

Table 5.2 Cost comparisons of Newton-Lanczos and Newton's Method  
using Example III ( contact problem ).

## APPENDIX A

## Orthogonalization of the Krylov Vectors

Consider the set of vectors  $(r_0, Ar_0, A^2r_0, \dots, A^{j-1}r_0)$ . Assume for the moment that a set of orthonormal vectors,  $(q_1, q_2, \dots, q_j)$ , are produced by successively orthogonalizing  $A^i r_0$  to all the previous vectors, all vectors  $A^k r_0$  for  $k < i$ , and then normalizing the resulting vector to obtain the next orthonormal vector  $q_{j+1}$ . Repeating this procedure for  $i = 1, 2, \dots, j-1$  defines the Gram-Schmidt method.

At some step  $i \leq j-1$  we have

$$\bar{r}_i = A^i r_0 - \sum_{k=1}^i \gamma_k q_k \quad (A1)$$

where  $\gamma_k = q_k^T A^i r_0$  and  $q_{i+1} = \frac{\bar{r}_i}{\|\bar{r}_i\|}$  but  $A^{i-1} r_0 \in \text{span}(q_1, q_2, \dots, q_i)$ . This can be rewritten as

$$A^{i-1} r_0 = \sum_{k=1}^i \nu_k q_k \quad (A2)$$

where  $\nu_k = q_k^T A^{i-1} r_0$ . Equations A1 and A2 can now be combined to give

$$\bar{r}_i = \sum_{k=1}^i \nu_k A q_k - \sum_{k=1}^i \gamma_k q_k \quad (A3)$$

and in turn  $A q_{i-1} \in \text{span}(q_1, q_2, \dots, q_i)$ . Therefore

$$\bar{r}_i = A q_i - \sum_{k=1}^i \bar{\gamma}_k q_k \quad (A4)$$

where  $\bar{\gamma}_k = q_k^T A q_i$ , and now  $q_{i+1} = \frac{\bar{r}_i}{\|\bar{r}_i\|}$ .

At some step  $j$ , the orthonormalization process described by equation A4 must be performed for the vector  $A q_j$ . However, consider the case when we orthogonalize this vector only to the previous two vectors,  $q_j$  and  $q_{j-1}$ . Then we get

$$r_j = Aq_j - \alpha_j q_j - \beta_j q_{j-1} \quad (A5)$$

where  $\alpha_j = q_j^T A q_j$  and  $\beta_j = q_{j-1}^T A q_j$ . The orthogonality of  $r_j$  to the remaining vectors can now be checked by forming  $q_i^T r_j$  for  $i < j - 1$

$$\begin{aligned} q_i r_j &= q_i^T (A q_j - \alpha_j q_j - \beta_j q_{j-1}) \\ &= q_i^T A q_j \\ &= q_j^T A q_i \end{aligned} \quad (A6)$$

This is derived using the orthogonality condition of the  $q$  vectors and the symmetry property of  $A$ . This equation together with equation A4 leads to the following:

$$q_i r_j = q_j^T (\bar{r}_i q_{i+1} + \sum_{k=1}^i \bar{r}_k q_k) = 0$$

This results shows that  $r_j$  is orthogonal to all the previous vectors  $q_i$ ,  $i = 1, 2, \dots, j$ . Therefore to generate the orthonormal bases, at each step  $Aq_j$  need only be orthogonalized against  $q_j$  and  $q_{j-1}$ .



## APPENDIX B

### Minimum Potential Energy over Krylov Space

The potential energy of a linear system is given by

$$E^n(x) = \frac{1}{2}x^T A x - x^T b \quad (B1)$$

where the  $E^n$  denotes the energy in the  $n$ -dimensional space. The minimum of this scalar function occurs when

$$\nabla E^n(x) = A x - b = 0 \quad (B2)$$

When the variables are restricted to the space of Lanczos vectors they are redefined as

$$x = Q_j f_j$$

and

$$b = \beta_1 q_1 = \beta_1 Q_j e_1$$

The potential energy now becomes

$$\begin{aligned} E^j(f_j) &= \frac{1}{2}f_j^T Q_j^T A Q_j f_j - f_j^T Q_j e_1 \beta_1 \\ &= \frac{1}{2}T_j f_j - f_j^T (\beta_1 e_1) \end{aligned} \quad (B3)$$

The minimum of the potential energy, restricted to the Krylov space, now occurs when

$$\nabla E^j(f_j) = T_j f_j - \beta_1 e_1 = 0 \quad (B4)$$

## References

- [1] K. W. Brodlie, A. R. Gourlay and J. Greenstadt, "Rank-One and Rank-Two Corrections to Positive Definite Matrices Expressed in Product Form," *J. Inst. Math. Appl.*, v. 11, pp. 73-82, 1973.
- [2] M. A. Crisfield, Letter to the Editor, *Int. J. Num. Meth. in Engr.*, v. 15, pp. 1419-1420, 1980.
- [3] J. R. Dennis, Jr. and J. J. More, "Quasi-Newton Methods, Motivation and Theory," *SIAM Review*, v. 19, pp. 46-86, 1977.
- [4] M. E. Gurtin, *Topics in Finite Elasticity*, CBMS-NSF Regional Conference series in Applied Mathematics, 35, SIAM 1981.
- [5] M. R. Hestenes and E. Stiefel, "Methods of Conjugate Gradients for Solving Linear Systems," *J. Res. Bur. Standards*, v. 49, pp. 409-436, 1952.
- [6] A. S. Householder, *The Theory of Matrices in Numerical Analysis*, Blaisdell, 1964.
- [7] S. Kaniel, "Estimates for some Computational Techniques in Linear Algebra," *Math. Comp.*, v. 20, pp. 369-378, 1966.
- [8] C. Lanczos, "An Iteration Method for the Solution of the Eigenvalue Problem of Linear Differential and Integral Operators," *J. Computational Phys.*, v. 26, pp. 43-65, 1977.
- [9] C. Lanczos, "Solution of Systems of Linear Equations by Minimized Iterations," *J. Res. Nat. Bur. Standards*, v. 49, pp. 33-53, 1952.
- [10] H. Matthies and G. Strang, "The Solution of Nonlinear Finite Element Equations," *Int. J. Num. Meth. Engr.*, v. 14, pp. 1613-1626, 1979.
- [11] L. Nazareth, "On the BFGS Method Part I: Properties of The Standard Algorithm," submitted for publications.
- [12] L. Nazareth, "On the BFGS Method Part II: Extended and Limited Memory Versions," submitted for publications.
- [13] N. M. Newmark, "A Method of Computation for Structural Dynamics," *J. Eng. Mech. Div., ASCE*, v. 85, EM3, pp. 67-94, 1959.
- [14] M. Ortiz, *Topics in Constitutive Theory for Inelastic Solids*, Ph.D. Thesis, University of California, Berkeley, 1981.
- [15] C. C. Paige, *The Computation of Eigenvalues and Eigenvectors of Very Large Sparse Matrices*, Ph.D. Thesis, University of London, 1971.

- [16] C. C. Paige, "Computational Variants of the Lanczos Method for the Eigenproblem," *J. Inst. Math. Appl.*, v. 10, pp. 373-381, 1972.
- [17] C. C. Paige and M. A. Saunders, "Solution of sparse indefinite systems of linear equations," *SIAM J. Num. Anal.*, v. 12, pp. 617-629, 1975.
- [18] B. N. Parlett, *The Symmetric Eigenvalue Problem*, Prentice-Hall, Englewood Cliffs (1980).
- [19] B. N. Parlett, "A New Look at the Lanczos Algorithm for Solving Symmetric Systems of Linear Equations," *Lin. Alg. Appl.*, v. 29, pp. 323-346, 1980.
- [20] B. N. Parlett and D. Scott, "The Lanczos Algorithm with Selective Orthogonalization," *Math. Comp.* v. 33, pp. 217-238, 1979.
- [21] J. C. Simo and R. L. Taylor, "Penalty Formulations for nonlinear incompressible elastostatics," to be published.
- [22] H. Simon, *The Lanczos Algorithm with Selective Reorthogonalization*, Ph.D. Thesis, University of California, Berkeley, 1982.
- [23] O. C. Zienkiewicz, *The Finite Element Method*, 3rd Edition, McGraw-Hill Book Co., London, 1977.